

## THESIS / THÈSE

### MASTER EN SCIENCES INFORMATIQUES

#### Apport du son, de la couleur et de la 3D à la représentation des Objets Symboliques

Baron, Xavier

*Award date:*  
1997

*Awarding institution:*  
Université de Namur

[Link to publication](#)

#### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal ?

#### Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

FACULTES  
UNIVERSITAIRES  
NOTRE-DAME DE LA PAIX  
NAMUR

---

INSTITUT D'INFORMATIQUE

**Apport du son, de la couleur  
et de la 3D à la représentation  
des Objets Symboliques.**

**Xavier Baron**

Promoteur :  
Professeur M. NOIRHOMME-FRAITURE

Mémoire présenté par  
**Xavier Baron**  
en vue de l'obtention du grade  
de Licencié et Maître en Informatique

---

**Année académique 1996-1997**

## Résumé

Les progrès technologiques *hardware*, en Bases de Données et dans les langages 'Orientés Objets' permettent de manipuler, stocker et représenter des objets exprimant des connaissances d'une complexité grandissante, difficiles à formaliser dans le carcan des tableaux de données classiques. C'est dans ce contexte que, sur base des travaux de DIDAY, nous introduisons la notion d'objets symboliques.

La perception de l'information est au coeur du problème de visualisation de données statistiques complexes. Cette complexité oblige le développement d'une méthode de représentation propre aux objets symboliques : « l'étoile zoom ».

Pour aider l'utilisateur à percevoir l'essentiel de l'information, nous nous intéressons au son, à la couleur et à la 3D comme supports de représentation, prémices de l'analyse visuelle d'objets. Notre objectif est d'éditer une série de recommandations pour utiliser au mieux ces trois moyens de présentation.

Plus en amont de ce problème de représentation visuelle, les traitements des données manquantes et de la sélection des variables les plus « utiles » à la représentation sont aussi étudiés.

## Abstract

Technological progress in hardware, in Data Bases and in Object Oriented languages is able to manipulate, to stock and to represent objects with data which are more and more difficult to be displayed. The complex data are not easy to be shown in classical tables. Based on DIDAY's idea, this problem leads to introduce the notion of symbolic objects.

The key of data visualization is the perception of information. We can cope with this difficulty by the development of a representation method used on symbolic object: « the Zoom Star ».

In order to help the user to visualize vital information, we will show some representation supports, early beginnings of visual perception such as the sound, the colour and the three dimensional forms. Our purpose is to condense a set of advices to use these three ways of presentation.

Furthermore, before the visualization aspect, missing data treatments and the selection of the most « useful » variables in representation aspect will be also examined.



## REMERCIEMENTS

---

C'est avec beaucoup de satisfactions que j'achève ce mémoire en exprimant ma profonde reconnaissance à tous ceux qui ont contribué à sa conception.

Elle s'adresse particulièrement à Madame M. **NOIRHOMME** qui m'en a proposé le sujet, pour sa persévérance malgré les « prolongations ».

Mes remerciements s'adressent aussi à Monsieur M. **ROUARD** pour les quelques réponses rapides par *e-mail*, aux questions de dernières minutes.

Je tiens à exprimer toute ma gratitude aux responsables de **la hiérarchie informatique des Assurances Fédérales** pour leur confiance témoignée à l'occasion de mon engagement, malgré la session et le mémoire inachevés.

Je remercie également mes **parents** et **amis** pour leur soutien moral et aide matérielle. Ils m'ont permis de mener à bien ce second cycle d'études universitaires et la réalisation de ce travail.

Une pensée particulière à l'adresse de mes **parents**, parents sollicités dans les moments les plus difficiles qu'impose un « bisseur » à son entourage.

Que tous ceux et celles qui m'ont donné un coup de main, si petit soit-il, et en particulier durant le *rush* final, puisent, dans ces quelques lignes, l'expression de mes plus vifs remerciements.

Bonne lecture.



## INTRODUCTION

---

L'évolution technologique en informatique est telle que les possibilités offertes aux utilisateurs augmentent considérablement de jour en jour.

*Comment utiliser ces nouvelles potentialités pour mieux percevoir, voire analyser le monde qui nous entoure ?*

De la caractérisation d'objets complexes du réel à leur représentation, une série de questions structurent notre mémoire. Au sein de chaque chapitre nous développons des éléments de réponse issus de recherches et d'expériences menées dans des domaines très variés.

*Comment caractériser les objets complexes, les objets symboliques ?*

L'arrivée de systèmes de saisie automatique d'informations dans tous les secteurs d'activité provoque une accumulation de données. Extraire de celles-ci des informations utiles est un problème majeur dans les entreprises. Grâce aux facilités de calculs et de représentations fournies par les ordinateurs de plus en plus performants, l'analyse des données a pour objet de résoudre partiellement cette situation.

Les progrès en Bases de Données et dans les langages 'Orientés Objets' permettent de manipuler et de représenter des objets exprimant des connaissances d'une complexité grandissante, impossible à exprimer dans le carcan des tableaux classiques de données.

Dans ce contexte et sur base des travaux de DIDAY, nous introduisons au **chapitre 01**, la notion d'objets symboliques.

L'aspect descriptif des données ainsi présentées permet déjà d'exprimer nombre d'informations et d'en tirer des enseignements. Au-delà de cette formalisation, nous voulons souligner les possibilités d'analyse que permettent les objets symboliques. Après avoir situé l'analyse des objets symboliques par rapport à l'analyse statistique classique et son intérêt dans différentes disciplines, nous détaillons dans le **chapitre 02**, deux études représentatives des potentialités de l'analyse symbolique.

En fin de chapitre, nous présentons un indice de dissimilarité permettant de comparer les objets.

*Comment représenter les objets symboliques ?*

Au-delà des traitements automatiques qui vont pouvoir être effectués sur les objets, il est « humainement » très intéressant de pouvoir s'en créer une « image ».

La technique habituellement utilisée est la représentation graphique des données, c'est à dire le passage d'une présentation chiffrée ou codée d'un individu ou d'un objet, à une représentation schématique.

Une représentation graphique efficace est une image synthétique, et donc nécessairement partielle, de l'individu. Le graphique et le tableau de nombres sont complémentaires. Le graphique fournit un cliché, une synthèse aisément interprétable du



tableau complet de nombres souvent peu lisible. Le graphique permet une meilleure compréhension et souvent une meilleure comparaison visuelle entre individus.

La perception de l'information est au coeur du problème de visualisation de données statistiques complexes. Cette complexité oblige le développement de méthodes de représentation propres aux objets symboliques. La méthode « l'étoile zoom », représentation graphique proposée par M. NOIRHOMME et M. ROUARD est présentée au **chapitre 03**.

*Au-delà de la représentation en étoile, comment aider l'utilisateur à percevoir l'essentiel de l'information ?*

Après avoir évoqué l'interfaçage au **chapitre 04**, nous nous intéressons au son, à la couleur et à la 3D comme supports de représentation, prémices de l'analyse visuelle d'objets.

Notre but est d'éditer une série de recommandations pour utiliser au mieux ces trois moyens de présentation. Les conclusions et la synthèse des chapitres les résumeront.

L'emploi du son, vocal ou non, au sein d'applications très interactives, est devenu de plus en plus courant et populaire en raison de son énorme potentiel. Il est utilisé pour présenter de l'information qui ne pourrait être l'objet de représentation via un mode de visualisation classique ou d'informations difficiles de discerner. Le son est habituellement un complément aux *outputs* visuels. Il augmente la quantité d'informations communiquées aux utilisateurs et/ou réduit la quantité d'informations que l'utilisateur doit traiter en mode visuel.

Au sein du **chapitre 05**, nous reprenons les résultats de nombreuses expériences montrant le rôle important que le son joue dans l'amélioration de la perception de l'information. Pour les personnes malvoyantes ou aveugles, utiliser le son représente parfois l'un des moyens accessibles de perception. Des recherches très poussées et des résultats concrets permettent d'identifier l'intérêt et les réelles potentialités qu'offre le son comme moyen de perception.

A sa guise, chaque utilisateur peut insérer des sons et même des séquences vocales enregistrées. Des utilitaires, comme ceux des paramètres du panneau de configuration de WINDOWS 95, permettent d'insérer des sons enregistrés pour une longue série d'événements pré-définis. Nous les appelons sons de signalisation. Nous développerons peu le sujet, si ce n'est que des expériences montrant leur efficacité. Nous décrivons plus en profondeur les sons comme moyens de perception et d'aide à la compréhension de phénomènes.

Dans le **chapitre 06**, au départ de l'expérience de TAPP, nous relatons l'importance de l'usage de la couleur. Nous décrivons certains phénomènes chromatiques et visuels ayant incidence dans la coloration des images symboliques. Sur base des disponibilités informatiques, nous proposons des teintes de référence pour l'application aux éléments des images des étoiles zoom.

Nous insistons dans le **chapitre 07** sur les raisons de l'utilisation de la 3D. Lors de la représentation 3D, certains éléments se retrouvent cachés par d'autres. Nous abordons ces situations de visibilité réduite en proposant des voies de contournements. Des outils particuliers permettant la manipulation des représentations 3D sont analysés.



Plus en amont de cette représentation visuelle, le traitement des données manquantes et la sélection des variables les plus « utiles » à la représentation sont également abordés.

En effet, l'un des buts de la représentation d'objets symboliques étant de les comparer, l'apparition de « trous » dans le tableau de nombres de base pose un sérieux problème.

D'autres part la représentation graphique donne une vue d'ensemble, le nombre de variables représentables sur celui-ci est forcément limité.

*Comment sélectionner les variables qui discriminent le mieux les objets symboliques entre eux ? Comment sélectionner les variables dignes d'intérêt pour la représentation graphique ?*

Il est possible de représenter graphiquement les objets symboliques. Représenter des objets symboliques à  $p$  variables peut ne pas être réalisable si  $p$  est très grand. De plus, cette représentation à  $p$  variables peut s'avérer non judicieuse : des variables sont fortement corrélées et leur représentation encombre le graphique, réduisant sa clarté et son aspect analytique.

Dans le **chapitre 08**, nous présentons des critères, des méthodes de sélection des variables. Identifier les variables qui, ensemble, permettent le mieux de discriminer les individus d'une population : tel est l'objectif à atteindre.

Dans le cas de la représentation graphique d'objets symboliques, en se fixant à 16 le maximum d'axes représentables sur une étoile zoom, on choisira en fonction d'un critère ou d'une méthode, le meilleur modèle contenant au plus 16 variables.

*Comment contourner les pertes d'informations causées par les données manquantes ?*

Lors d'enquêtes, des individus interrogés ne répondent pas ou de manière incohérente ou partielle. Lors de prises automatiques continues de mesures, un appareil défectueux biaise certaines données. Ces données manquantes ne doivent pas altérer les traitements des analyses, il faut, soit les éliminer, soit les remplacer.

Dans le **chapitre 09**, nous nous intéressons à différentes méthodes de résolution de ce type de situations.

En toute généralité, les données manquantes impliquent directement une augmentation de l'imprécision dans l'analyse des résultats. Cette imprécision dépend du nombre de données manquantes. Il faut en tenir compte dans l'interprétation des résultats.



### 1. Introduction

Ce premier chapitre a pour principal objectif d'introduire les différentes notions de l'analyse des données symboliques. Le symbolisme adopté correspond à celui des objets symboliques introduits par E. DIDAY en 1987 [DIDAY 87]. Nous présentons ici un éventail de définitions et d'exemples de base. Pour une présentation complète, se référer [DIDAY 91].

Les références bibliographiques utilisées pour la conception de ce chapitre sont les suivantes : [DIDAY 87], [DIDAY 91], [DIDAY 93], [DE CARVALHO 94], [LÊ 96], [PERINEL 92] et [REGNIER 92].

### 2. Les concepts de la théorie des objets symboliques

L'analyse des données symboliques s'assigne plusieurs objectifs :

- savoir représenter nos connaissances par des expressions à la fois symboliques et numériques;
- savoir manipuler et utiliser ces expressions dans le but d'aider à décider, de mieux analyser, synthétiser et organiser les expériences et les observations.

Les objets symboliques présentés dans ce chapitre constituent les individus de l'analyse des données symboliques. Ils permettent de représenter des individus complexes ou des classes d'individus par des conjonctions de propriétés. Certaines variables peuvent prendre des valeurs multiples et pondérées et sont parfois reliées entre elles par des relations d'ordre logique. Des outils pour manipuler ces objets sont décrits : union, intersection, généralisation, extension, ...

#### 2.1. Les variables

Une variable  $y$  est par définition une de fonction de  $\Omega \rightarrow O$  où :

- $\Omega$  est l'ensemble des objets élémentaires et,
- $O$  l'ensemble des observations où la variable  $y$  prend ses valeurs.

Selon la structure algébrique de  $O$ , le type de la variable est différent. La variable est dite quantitative si  $O$  est un intervalle de l'ensemble des réels. Si  $O$  est fini ou dénombrable, nous dirons que la variable est qualitative. Dans ce cas la variable est qualitative ordinale ou nominale selon que  $O$  est ordonné ou non.

Exemple :

Soit  $\Omega = \{\text{Berlingo, Courrier, Combo, Expert, Express}\}$  et les variables : la classe de taxe de mise en circulation, la nationalité et la cylindrée du moteur définissant le tableau 1.1. suivant. Ces variables notées  $y_1, y_2, y_3$  sont respectivement qualitatives ordinale, nominale et quantitatives.



	y <sub>1</sub>	y <sub>2</sub>	y <sub>3</sub>
Berlingo	2	France	1905
Courier	2	Allemagne	1753
Combo	1	Allemagne	1389
Expert	2	France	1905
Express	2	France	1870

Tableau 1.1. : exemple de données et types de variable.

## 2.2. Les objets symboliques

Une ligne d'un tableau de données caractérise un individu. Elle peut s'exprimer sous forme d'une conjonction de propositions logiques appelées événements. Si l'on reprend l'exemple du tableau 1.1., le premier individu ou l'objet Berlingo, il est caractérisé par l'expression symbolique définie par la conjonction :

$[ \text{classe de taxe de mise en circulation} = 2 ] \wedge [ \text{nationalité} = \text{France} ] \wedge [ \text{cylindrée} = 1905 ]$ .

On définit un objet symbolique comme *une description qui s'exprime à l'aide d'une conjonction événements (on dit aussi propriétés) portant sur des valeurs prises par les variables*.

Pour avancer graduellement dans la compréhension, nous définissons d'abord les objets symboliques dans le cas où ce sont des fonctions à valeurs dans {vrai, faux} à une variable.

Ensuite nous introduisons les objets symboliques munis de propriétés où de telles fonctions à plusieurs variables peuvent apparaître.

Trois raisons justifient l'utilisation du terme « objet » :

- 1 - faisant référence à l'analyse classique des données, ce terme permet de rappeler que les objets symboliques peuvent être considérés comme des lignes d'un tableau de données sur lesquelles on tentera d'adopter les méthodes usuelles d'analyse des données,
- 2 - les objets symboliques tels qu'ils sont présentés peuvent être considérés comme des « objets » au sens de la programmation Orientée Objet,
- 3 - les événements de la connaissance supplémentaire ne seront pas qualifiés d'objets, même s'ils peuvent prendre la même forme que les objets analysés.

Dans les langages de programmation Orientés Objets, chaque objet comprend une partie purement déclarative qui décrit les variables et une partie active formée des méthodes particulières pour calculer certaines variables ou valeurs. Les objets sont organisés selon un graphe dit « d'héritage » à partir d'un objet générique, père de tous les autres. De cette façon, chaque objet apparaît comme une instanciation d'un ou plusieurs objets qui affinent la description. Le parallélisme entre les objets symboliques et les langages orientés objets existe aussi bien dans leur définition que dans leur organisation (obtenue par classification).



Exprimé sous sa forme logique par l'expression symbolique qui le représente, un objet symbolique  $s$  est dit défini en intention ou en compréhension. L'ensemble des événements élémentaires de  $\Omega$  pour lesquels il est satisfait (vrai) constitue son extension notée  $|s|_\Omega$ .

### 2.3. Les événements élémentaires

Soit  $p$  variables  $y_1, \dots, y_p$  définies sur  $\Omega$ , l'ensemble des objets élémentaires observés et prenant leurs valeurs dans  $O_1, \dots, O_p$ . Si l'on identifie chaque élément de  $\Omega$  avec l'ensemble des valeurs qu'il prend, on peut plonger  $\Omega$  dans  $\Omega' = O_1 \times \dots \times O_p$  ensemble des objets élémentaires possibles. On note  $V_i$  une partie de  $O_i$ .

Définition :

Un événement élémentaire représenté par l'expression symbolique  $e_i = [y_i = V_i]$  est défini par la fonction  $e_{y_i V_i}$  :

$$\Omega \rightarrow \{\text{vrai, faux}\} \text{ telle que } e_{y_i V_i}(w) = \text{vrai si et seulement si } y_i(w) \in V_i.$$

Ainsi dans le cas d'un objet symbolique défini par un événement élémentaire  $[y_i = V_i]$ , sa définition en compréhension notée  $e$  est  $e = [y_i = V_i]$ .

Son extension est :

$$|e|_\Omega = \{w \in \Omega / y_i(w) \in V_i\}, \text{ où } w \text{ représente le nom d'un objet symbolique.}$$

Pour simplifier les notations, et comme le propose l'auteur [DIDAY 91], on considère que  $e_i$  exprime aussi bien le nom de l'événement élémentaire que celui de son expression symbolique  $[y_i = V_i]$ .

Exemple :

Considérons la troisième variable du tableau 1.1., l'événement élémentaire noté  $e = [y_3 = \{1870, 1905\}]$  est défini par la fonction  $e_{y_3 V_3} : \Omega \rightarrow \{\text{vrai, faux}\}$  telle que  $e_{y_3 V_3}(w) = \text{vrai si et seulement si } y_3(w) \in \{1870, 1905\} = V_3$ .

Son extension est  $|e|_\Omega = \{\text{Berlingo, Expert, Express}\}$ .

### 2.4. Les objets assertion.

Un objet assertion est une conjonction d'événements élémentaires.

Définition :

Un objet assertion de représentation symbolique  $a = [y'_1 = v_1] \wedge \dots \wedge [y'_q = v_q]$  où  $V_i \subset O'_i$  est défini par la fonction  $a_{y' V} : \Omega \rightarrow \{\text{vrai, faux}\}$  telle que  $a_{y' V}(w) = \text{vrai si et seulement si pour tout } i = 1, \dots, q \text{ on a } y_i(w) \in V_i$ .

Si  $a$  est une conjonction d'événements élémentaires qui doivent être vrais simultanément pour le même objet élémentaire  $w$ , on note  $a$  sous la forme symbolique suivante :

$$a = [y'_1(w) = V_1] \wedge \dots \wedge [y'_q(w) = V_q].$$

Cette notation est réservée aux objets hordes vus plus loin.



L'ensemble des objets élémentaires qui satisfont l'objet assertion  $a$  dans  $\Omega$  est noté  $|a|_\Omega$  et constitue la définition en extension de  $a$  dans  $\Omega$  :

$$|a|_\Omega = \{w \in \Omega / y'_i(w) \in V_i, \text{ pour } i = 1, \dots, q\}$$

Exemple :

Reprenons le  $\Omega$  présenté dans le tableau 1.1. Un objet élémentaire tel que  $y_1 = 1$  et  $y_2 = \text{Allemagne}$  s'exprime sous la forme de l'objet assertion suivant :

$$a = [y_1 = 2] \wedge [y_2 = \text{Allemagne}].$$

L'ensemble des objets élémentaires de  $\Omega$  dont la nationalité d'origine est soit France, soit Allemagne et la cylindrée moteur comprise entre 1800 et 1910 peut s'écrire :

$$a = [y_2 = \{\text{France, Allemagne}\}] \wedge [y_3 = [1800, 1910]]$$

Pour chaque  $w_i \in \Omega$  (élément de  $\Omega$ ), on peut associer une assertion  $a_i$  telle que :

$$\begin{aligned} a_{\text{Berlingo}} &= [y_1 = 2] \wedge [y_2 = \text{France}] \wedge [y_3 = 1905] \\ a_{\text{Courier}} &= [y_1 = 2] \wedge [y_2 = \text{Allemagne}] \wedge [y_3 = 1753] \\ a_{\text{Combo}} &= [y_1 = 1] \wedge [y_2 = \text{Allemagne}] \wedge [y_3 = 1389] \\ a_{\text{Expert}} &= [y_1 = 2] \wedge [y_2 = \text{France}] \wedge [y_3 = 1905] \\ a_{\text{Express}} &= [y_1 = 2] \wedge [y_2 = \text{France}] \wedge [y_3 = 1870] \end{aligned}$$

Définissons l'objet assertion :

$$a = [y_1 = \{2,3\}] \wedge [y_2 = \{\text{Espagne, France}\}] \wedge [y_3 = [1700, 2000]];$$

dont nous allons chercher l'extension :

$$\begin{aligned} a &= [y_1 = \{2,3\}] = [y_1 = 2] \vee [y_1 = 3] \\ &\wedge [y_2 = \{\text{Espagne, France}\}] = [y_2 = \text{Espagne}] \vee [y_2 = \text{France}] \\ &\wedge [y_3 = [1700, 2000]]; \end{aligned}$$

Par développement, nous obtenons :

$$\begin{aligned} a &= [y_1 = 2] \wedge [y_2 = \text{Espagne}] \wedge [y_3 = [1700, 2000]] \\ &\vee [y_1 = 2] \wedge [y_2 = \text{France}] \wedge [y_3 = [1700, 2000]] \\ &\vee [y_1 = 3] \wedge [y_2 = \text{Espagne}] \wedge [y_3 = [1700, 2000]] \\ &\vee [y_1 = 3] \wedge [y_2 = \text{France}] \wedge [y_3 = [1700, 2000]] \end{aligned}$$

Les seuls objets de  $\Omega$  qui satisfont  $a$  sont : Berlingo, Expert, Express; autrement dit, l'extension de  $a$  dans  $\Omega$  est :

$$|a|_\Omega = \{ \text{Berlingo, Expert, Express} \}.$$

## 2.5. Les objets hordes.

Reprenons le tableau 1.1. et considérons les événements élémentaires  $e_1 = [y_1 = 1]$  et  $e_2 = [y_3 = 1753]$ ;  $e_1(\text{Combo}) = \text{vrai}$  et  $e_2(\text{Courier}) = \text{vrai}$ . Par contre il n'existe pas d'individu de  $\Omega$  pour lesquels  $e_1$  et  $e_2$  soient vrais simultanément.

Définition :

Un objet  $h$  représenté par l'expression symbolique  $h = [y'_1(u_1) = V_1] \wedge \dots \wedge [y'_p(u_p) = V_p]$  est défini par la fonction  $h_{yV} : \Omega^q \rightarrow \{\text{vrai, faux}\}$  telle que  $\forall W = (w'_1, \dots, w'_q) \in \Omega^q$ ,  $h_{yV}(W)$  est vrai si et seulement si  $\forall i \ y'_i(w'_i) \in V_i$ .

Le objet assertion constitue un cas particulier d'objet horde.

L'extension d'un objet horde : soit  $h$  un objet horde défini sur  $\Omega^q$ . L'extension de  $h$  est l'ensemble des éléments  $W \in \Omega^q$  tel que  $h_{yV}(W) = \text{vrai}$ . On note cette extension  $|h|_\Omega$ .

Considérons dans l'exemple du tableau 1.1., l'objet horde suivant :

$$h = [y_1(u_1) = 2] \wedge [y_2(u_2) = 1905],$$

on a  $|h|_\Omega = \{(\text{Berlingo}, \text{Berlingo}), (\text{Berlingo}, \text{Expert}), (\text{Courier}, \text{Berlingo}), (\text{Courier}, \text{Expert}), (\text{Expert}, \text{Expert}), (\text{Express}, \text{Berlingo}), (\text{Express}, \text{Expert})\}$ .

## 2.6. Les objets de synthèse.

Définition :

Un objet de synthèse est la conjonction de  $k$  objets hordes respectivement définis sur chacun des ensembles  $H_1, \dots, H_k$ . Il s'écrit sous la forme générale  $s = h_1 \wedge \dots \wedge h_k$  avec  $h_i \in H_i$ .

Exemple :

A partir des murs ( $\Omega_0$ ), des fenêtres ( $\Omega_1$ ), des portes ( $\Omega_2$ ) et des toits ( $\Omega_3$ ) on définit des objets de synthèse appelés « Type de maison ».

## 2.7. Les objets symboliques munis de méthodes et de propriétés.

On généralise les différents objets définis jusqu'à présent à des objets munis de méthodes, de propriétés. On ajoute par conjonction des événements définissant des méthodes (par exemple : calculer une variable ou une fonction) ou des propriétés exprimant des liens entre objets, entre variables, entre variables et objets.

Notons  $\text{Meth}(y_1, \dots, y_k)$  la méthode de calcul de la variable  $y_i$  qui dépend des variables  $y_1$  à  $y_k$  (par exemple :  $y_i$  est une surface et la méthode : calcul de cette surface en fonction des dimensions des côtés).

Exemple d'objets assertion munis de méthodes :

$\text{Voiture}_p = [\text{consommation moyenne} = [4, 10]] \wedge [\text{prix carburant} = [23, 40]] \wedge [\text{nombre kilomètres annuels} = \{5000, 10000, 15000, 20000\}] \wedge [\text{coût annuel} = f(c, n, p)]$

Où  $f(c, n, p)$  est une méthode permettant le calcul du coût annuel d'une voiture (qui n'est pas une variable) à partir de la consommation moyenne ( $c$ ), du nombre de kilomètres parcourus annuellement ( $n$ ) et du prix du carburant ( $p$ ).

La propriété peut également s'exprimer sous forme de règle. Par exemple : signaler que si la variable  $y_1$  prend la valeur  $a$  alors la variable  $y_2$  prend obligatoirement la valeur  $c$ ; si la variable  $y_1$  prend la valeur  $b$ , la variable  $y_2$  prend obligatoirement la valeur  $d$ .



On écrit :

$$y_1 = \{a,b\} \quad \wedge \quad y_2 = \{c,d\} \\ \wedge \quad [[\text{si } [y_1 = a] \text{ alors } [y_2 = c]] \vee [\text{si } [y_1 = b] \text{ alors } [y_2 = d]]].$$

Ou sous la forme :

$$y_1 = \{a,b\} \quad \wedge \quad y_2 = \{c,d\} \\ \wedge \quad [\text{si } [y_1 = a] \text{ alors } [y_2 = c]] \\ \wedge \quad [\text{si } [y_1 = b] \text{ alors } [y_2 = d]]$$

Les notions de méthodes et de propriétés peuvent être généralisées aux objets hordes et de synthèse. Dans le cas des objets hordes, les méthodes peuvent dépendre de plusieurs variables et/ou objets élémentaires. Les objets de synthèse peuvent également être munis de méthodes : dans l'exemple de la maison décrite par les objets élémentaires murs, fenêtres, portes et toits, une méthode situe les fenêtres par rapport aux portes et une autre exprime la taille d'une fenêtre à celle d'une porte.

## 2.8. La connaissance supplémentaire.

Il faut distinguer les objets symboliques que l'on étudie des connaissances supplémentaires que l'on calcule sur ces objets. A titre d'exemples de connaissances supplémentaires : des mesures de ressemblance, des dénominations regroupant des modalités, des contraintes exprimées sous forme de règles. Avec l'analyse symbolique : les mesures de ressemblance et les contraintes peuvent s'exprimer par des objets assertions, hordes ou règles.

### 2.8.1. Les taxonomies sur les objets élémentaires.

Au point 2.1. de ce chapitre, nous définissions une variable  $y_i$  comme une fonction de  $\Omega$ , ensemble des objets élémentaires dans  $O_i$  : un ensemble d'observations.  $y_i$  est une application sur une partie de  $\Omega$ ; chaque  $y^{-1,i}$  définit une partition sur cette partie. Si  $y^{-1,i}$  est un recouvrement et non plus une partition, on dit que  $y^{-1,i}$  est une variable taxonomique.

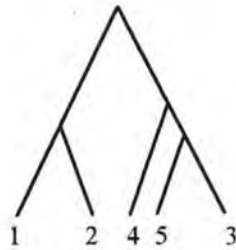
Les variables hiérarchiques et pyramidales sont deux types de partitions de variables taxonomiques.

Une variable hiérarchique est une correspondance  $t_h$  de  $\Omega \rightarrow O$ , l'ensemble des parties de  $\Omega$  associant à tout élément de  $O$  l'ensemble des paliers de la hiérarchie auxquels cet élément appartient. Si l'on note  $t_h^{-1}(\sigma)$  l'ensemble des  $w \in \Omega : \sigma \in t_h(w)$  on a  $\forall \sigma, \sigma' \in O$  :

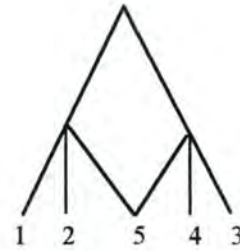
$$t_h^{-1}(\sigma) \cap t_h^{-1}(\sigma') = \begin{cases} \emptyset \\ \text{Ou } t_h^{-1}(\sigma) \subset t_h^{-1}(\sigma') \\ \text{Ou } t_h^{-1}(\sigma') \subset t_h^{-1}(\sigma) \end{cases}$$



Une variable  $t_p$  est pyramidale si  $\forall \sigma, \sigma' \in O : t_p^{-1}(\sigma) \cap t_p^{-1}(\sigma')$  est vide ou s'il existe  $\sigma''$  tel que  $t_p^{-1}(\sigma) \cap t_p^{-1}(\sigma') = t_p^{-1}(\sigma'')$ .



Hierarchie



Pyramide

La hiérarchie est une suite de partition emboîtées de moins en moins fines dont la représentation est l'arbre hiérarchique ou dendrogramme. Le modèle hiérarchique est étendu au modèle pyramidal qui construit une suite de recouvrements emboîtés au lieu de partitions emboîtées.

## 2.9. Le traitement des données symboliques.

Les données classiques se présentent sous forme d'un ensemble d'objets élémentaires caractérisés par un nombre fini  $p$  de variables qualitatives et quantitatives. On considère ces objets comme des points de l'espace  $\mathbb{R}^p$  en transformant les variables qualitatives en variables binaires. Ces données sont qualifiées de numériques.

En analyse des données, la connaissance supplémentaire est définie par le choix d'une mesure de dissimilarité entre les objets. La mesure est utilisée pour réaliser une classification des objets ou pour réaliser une représentation plane par analyse factorielle.

Les objets symboliques se présentent sous forme de conjonctions d'événements. Ils sont de types différents : sous forme d'objets élémentaires, d'objets assertion, d'objets hordes, d'objets de synthèse ou d'objets règles (tous avec ou sans méthodes et propriétés). La connaissance supplémentaire est définie par les taxonomies, les affinités et les règles. Elle peut concerner des objets de même type ou se référer à des combinaisons plus complexes. Les objets de synthèse « Types de maison » peuvent être considérés comme des combinaisons de murs, de portes et de fenêtres, eux aussi des objets de synthèse.

### 2.9.1. Le traitement des données et les quatre principes.

Comme pour l'analyse des données classiques (numériques) l'objectif est de résumer, et de synthétiser l'information contenue dans les données symboliques pour constituer une base à un processus décisionnel, de reconnaissance ou plus généralement dans un but descriptif.

Le premier point de vue par lequel on aborde l'analyse est du type numérique : dans la mesure du possible, on adapte les méthodes classiques de l'analyse des données aux nouveaux objets.



Le second point de vue relève du monde symbolique : il consiste à adopter les quatre principes suivants :

- 1) *Le principe de fidélité* : les données doivent être fidèles à la réalité multidimensionnelle qu'elles tentent de représenter, en évitant les artefacts des codages réducteurs. Par exemple : on ne réduit pas à une moyenne les valeurs multiples présentes.
- 2) *Le principe de la prédominance de la connaissance* : la connaissance supplémentaire, qu'elle soit fournie par l'interactivité homme-machine ou machine-machine, guide les algorithmes. Pour expliciter ce principe, citons une méthode de l'analyse des données classiques où les algorithmes sont guidés par les questions posées : la segmentation, où à chaque réponse le système propose une série de nouvelles questions par ordre de pouvoir discriminant.  
L'approche symbolique est propice à l'interaction homme-machine, les résultats sont aisément interprétables. Un processus d'apprentissage automatique peut fournir des exemples et des contre-exemples par simulation afin de tester la qualité des règles obtenues.
- 3) *Le principe de cohérence* : les résultats des analyses devront s'exprimer selon les mêmes termes que les objets fournis au départ. Ce principe est appliqué de façon implicite en analyse classique. L'application de ce principe dans le domaine symbolique entraîne que l'expression des résultats en termes d'objets assertion, hordes, règles ou de synthèse.
- 4) *Le principe d'explicabilité* : fournir des résultats explicables et d'utilisation aisée sont nécessaires pour définir une base de connaissances dans le domaine d'où sont issues les données.

## **2.9.2. L'apprentissage de la connaissance.**

Au départ d'un ensemble d'objets A (pouvant être dynamique : aller en s'agrandissant) on cherche à obtenir un ensemble plus simple d'objets B. Cet ensemble peut aussi être dynamique (aller en s'améliorant). Les ensembles A et B sont chacun composés d'objets d'un même type ou non; les types des deux ensembles peuvent être différents.

## **2.10. Propriétés des objets symboliques.**

Le lecteur intéressé trouvera dans la littérature spécialisée [DIDAY 87], [DIDAY 91], [DIDAY 93], [DE CARVALHO 94] et [PERINEL 92] un ensemble de propriétés et de leurs démonstrations. Nous reprenons ici les plus courantes : celles utilisées dans les analyses décrites au chapitre suivant.

*L'ordre symbolique :*

$\forall s_1, s_2 \in S$ , on dit que  $s_1 \leq s_2$  si et seulement si  $s'_1 \subseteq s'_2$ .



où :

$s_1$  et  $s_2$  sont des objets symboliques et  $s'_1, s'_2$  les extensions respectives. L'ordre est une relation partielle réflexive et transitive. Cette relation n'est pas antisymétrique puisque deux objets de même extension ne sont pas forcément identiques.

*L'héritage et la généralisation :*

$\forall s_1, s_2 \in S$ , on dit que  $s_1$  hérite de  $s_2$  et que  $s_2$  est plus général que  $s_1$  si et seulement si  $s'_1 \subset s'_2$ . L'ensemble des objets symboliques  $S$  tels que  $s (\in S) \geq s_1$  sont les ascendants de  $s_1$ . L'ensemble des objets symboliques  $s$  tels que  $s \leq s_1$  sont les descendants de  $s$ .

*L'union et l'intersection symbolique :*

L'union symbolique  $s_1 \cup s_2$  est la conjonction de tous les objets symboliques de  $S$  dont l'extension symbolique contient l'ensemble des objets symboliques de  $s'_1$  et  $s'_2$ . L'intersection symbolique  $s_1 \cap s_2$  est la conjonction de tous les objets symboliques de  $S$  dont l'extension symbolique contient l'ensemble des objets symboliques communs à  $s'_1$  et  $s'_2$ .

*La I-extension :*

Dans le cas où l'on ne dispose pas de  $\Omega$  mais seulement d'un ensemble  $S_1 (\subset S)$  d'objets symboliques, on définit la I-extension d'un objet symbolique  $s$  dans  $S_1$  comme étant l'ensemble des  $s_1 \in S_1$  tels que  $|s_1|_{\Omega} \subset I|s|_{\Omega}$ .

## 2.11. Qualité des objets symboliques.

Si j'observe des voitures, toutes d'origine allemande, et que je dis : « je vois des voitures », mon assertion ne décrit pas de façon complète mon observation puisque j'ometts de dire qu'elles sont d'origine allemande. Cette idée se traduit par la notion de complétude d'un objet symbolique. Formellement, cela revient à mesurer l'écart entre l'ensemble des événements élémentaires dont la conjonction définit un objet symbolique  $s$  et l'ensemble de tous les événements élémentaires dont l'extension contient  $s'$ .

Notations :

- $d(s)$  (d comme définition) est l'ensemble des événements élémentaires dont la conjonction définit  $s$ ,
- $c(s)$  (c comme complet) est l'ensemble de tous les événements élémentaires de plus petite extension symbolique qui contiennent  $s'$ ,
- $s^c$  est l'objet symbolique défini par la conjonction de tous les éléments de  $c(s)$ .

Exemple :

	$Y_1$	$Y_2$	$Y_3$
$w_1$	1	0	0
$w_2$	1	1	1
$w_3$	1	1	1
$w_4$	0	1	1



$$\begin{aligned}
s &= [y_1 = 1] \wedge [y_3 = 1] \\
d(s) &= \{[y_1 = 1], [y_3 = 1]\} \\
c(s) &= \{[y_1 = 1], [y_2 = 1], [y_3 = 1]\} \\
s^c &= [y_1 = 1] \wedge [y_2 = 1] \wedge [y_3 = 1]
\end{aligned}$$

### 2.11.1. Complétude d'un objet symbolique.

Voici une des manières de calculer la complétude d'un objet symbolique :

$$c_1(s) = \text{card}(c(s) \setminus d(s)).$$

En reprenant les données de l'exemple ci-avant : la complétude  $c_1(s) = 1$ .

Un objet est dit complet si et seulement si  $c(s) = d(s)$ ; autrement dit, si  $c(s) = 0$ . DIDAY dans [DIDAY 87] et [DIDAY 91] énonce et démontre des propriétés portant sur les objets complets. Le lecteur intéressé se référera à la bibliographie.

### 2.11.2. Affinement d'un objet symbolique.

On dit qu'un objet symbolique est d'autant plus affiné que les événements élémentaires qui le définissent ont une extension proche de celle de  $s$ . Le degré d'affinement d'un objet symbolique s'exprime via le critère suivant :

$$A(s) = \frac{\text{card}(\cup e'_i(s) \setminus \cap e'_i(s))}{\text{card}(\cup e'_i(s))}$$

où les  $e_i(s)$  sont les événements élémentaires de  $s$ .

Exemple :

En reprenant les données ci-dessus :

$$\begin{aligned}
s &= [y_1 = 1] \wedge [y_3 = 1] \\
A(s) &= \frac{1}{5} \text{card}(\Omega \setminus \{w_2, w_3\}) = \frac{2}{4} = \frac{1}{2}
\end{aligned}$$

### 2.11.3. Simplicité d'un objet symbolique.

On dit qu'un objet symbolique  $s$  est d'autant plus simple que le nombre d'événements élémentaires qui le décrit est proche d'un ensemble d'événements élémentaires de cardinal minimum. La conjonction de ces événements élémentaires est notée  $s_p$  et a la même extension. La simplicité peut se mesurer par exemple, à l'aide d'un critère de type :

$$S(s) = \text{card } d(s) - \text{card}(d(s_p)).$$

Exemple :

En reprenant les données ci-dessus :

$$\begin{aligned} s &= [y_1 = 1] \wedge [y_2 = 0] \\ s_p &= [y_3 = 0] \\ S(s) &= 2 - 1 = 1 \end{aligned}$$

#### 2.11.4. La redondance d'un objet symbolique.

On peut combiner la simplicité et la complétude d'un objet symbolique en définissant sa redondance. La redondance d'un objet symbolique est l'écart entre le cardinal de  $s_p$  et celui de  $c(s)$ . Un objet symbolique à la fois complet et simple a sa redondance nulle.

#### 2.12. Qualité des classes d'objets symboliques.

L'examen des qualités d'une classe demande de définir ce que sont : l'extension, l'ordre, l'union et l'intersection de classe. Une façon simple de procéder consiste à associer un objet symbolique à chaque classe. Par exemple : l'union ou l'intersection des objets de la classe. Les différentes notions, propriétés et qualités des objets symboliques sont alors appliquées à ces objets. Une classe est qualifiée complète si l'objet symbolique est lui-même complet.

L'extension d'une classe est l'union ou l'intersection de l'extension des objets de la classe. Pour définir l'ordre entre classes, l'une est inférieure une autre si son extension est contenue dans celle de cette autre.

Des propriétés caractéristiques de la notion de classe peuvent se distinguer : la stabilité et l'effritement.

**La stabilité d'une classe** est la capacité de celle-ci à être représentée par l'objet symbolique de plus petite extension qui contient l'union des extensions des éléments de la classe.

$$st(C) = \text{card}(\bigcup c_i |_{\Omega} - \bigcup | c_i |_{\Omega}).$$

Où :

C est la classe et les  $c_i$ , les éléments de cette classe.

Une classe C est stable si  $st(C) = 0$ .

**L'effritement d'une classe** est le plus petit nombre d'objets symboliques  $a \in A \subset S$  dont la réunion des extensions est contenue dans l'extension des éléments d'une classe C tout en en s'en écartant le moins possible. Etant donnée une classe  $C \subset S$  d'éléments  $c_i$  et une classe  $A \subset S$  d'éléments  $a_i$ , un des critères de mesure de l'effritement se calcule par :

$$E_2(C) = \text{Min} \{ \text{card } A | \bigcup | a_i |_{\Omega} = \bigcup | c_i |_{\Omega} \}.$$

Où :

s est tel que :  $\bigcup | c_i |_{\Omega} = \bigcup | a_i |_{\Omega}$ . Une classe C est d'effritement minimum lorsque  $E_2 = 1$ .



## 2.14. Une extension à des objets symboliques modaux.

### 2.14.1. Introduction aux objets symboliques modaux.

Pour représenter la réalité multidimensionnelle, la logique booléenne ne suffit pas. Dans la pratique, on est amené à faire intervenir des jugements ou modes portant sur des événements élémentaires. Appel est fait à la logique modale aux travers de différentes sémantiques :

- la sémantique aléthique pour exprimer : « il se peut que » / « il est certain que »;
- la sémantique déontique : « il est permis de » / « il est interdit de »;
- la sémantique épistémique : « il ignore si » / « il croit que »;
- la sémantique temporelle : « par le passé » / « à l'avenir ».

Pour toutes ces sémantiques, nous avons supposé deux modes opposés, mais de nombreuses variantes de ces modes sont imaginables. Par exemple l'usage d'une sémantique de « convenance » exprimant les différentes conditions à remplir pour « convenir » à un emploi, les modes suivants sont utilisés :  $M_1$  = « il est nécessaire que »,  $M_2$  = « si possible »,  $M_3$  = « il est impossible »,  $M_4$  = « il n'est pas nécessaire »; possible exprime ici une préférence plutôt qu'une possibilité au sens de la théorie possibiliste.

Les modes peuvent concerner globalement toutes les valeurs prises par une variable dans une classe donnée et donc porter sur un événement élémentaire en s'écrivant sous la forme :  $a_x = \hat{1} M_i [y_i = V_i]$  où  $x$  est associé à la sémantique liée à la connaissance du domaine étudié. Il s'agit d'objets modaux de l'extérieur.

Dans le cas où les modes portent sur les valeurs prises par les variables, nous obtenons des objets modaux dits de l'intérieur et ils sont de la forme :  $a_x = \hat{1} [y_i = M_i V_i]$ .

Il existe au moins deux façons de définir l'extension dans  $\Omega$  des objets symboliques modaux. La première consiste à considérer que chaque élément  $\omega \in \Omega$  est plus ou moins dans l'extension selon que son poids ( $a_x(\omega)$ ) est plus ou moins grand. L'extension est définie par tous les couples  $(a_x(\omega), \omega)$ . La seconde façon de procéder consiste à déterminer un seuil  $\alpha \in a(\Omega)$  et à considérer la  $\alpha$ -extension de  $a_x$  dans  $\Omega$  notée  $|a_x|_{\Omega, \alpha}$  comme étant formée de tous les éléments  $\omega \in \Omega$  tels que  $a_x(\omega) \geq \alpha$ .

### 2.14.2. Assertions modales de l'extérieur.

Définition.

Supposons les hypothèses suivantes satisfaites :

- $\forall i \exists j : M_i = \neg M_j$  (table des négations);
- $M_i e = M_j e \Rightarrow M_i = M_j$  (où  $e$  est un événement élémentaire booléen);
- $M_i e = \neg M_j \neg e$  (axiome des logiques modales).



Etant donnés :

- deux ensembles de modes  $M = \{M_i\}$ ,  $m = \{m_i\}$  et  $L = \{L_i\}$ ;
- une fonction  $g : m \times M \rightarrow L$  satisfaisant à la condition  $g(m_i, \neg M_j) = g(\neg m_i, M_j)$ ;
- une fonction d'agrégation  $f : P(L) \rightarrow L$
- des applications  $z : \Omega^s \rightarrow m$  telles que si  $\omega^s = \hat{e}_i$

alors  $z_i(\omega^s) = m_i$  est le mode associé à l'événement élémentaire  $e_i$ .

On définit un objet modal de l'extérieur comme une application  $a_{YV}$  de  $\Omega$  dans  $L^*$  où  $Y = (y_1, \dots, y_p)$  et  $V = (V_1, \dots, V_p)$  notée  $a = \hat{I}_M M_i [y_i = V_i]$  telle que si  $\omega^s = \hat{e}_i$  où les  $e_i$  sont normalisés par rapport à  $a$ , à l'aide d'une table de négation;

alors,  $a(\omega) = \hat{I}_M g(m_i, M_i) = \hat{I}_M L_i = f(L_1, L_2, \dots, L_k)$  si,  $i = 1, \dots, k$ .

Exemple :

Dans une annonce d'offre d'emploi, une entreprise indique qu'elle recherche une personne dont l'âge est nécessairement compris entre 25 et 35 ans, dont le diplôme est du type  $A_1$  ou  $A_2$  et qui ne peut être étrangère. Cette annonce s'exprime par l'objet modal de l'extérieur où  $M_1$  = il est nécessaire,  $M_2$  = si possible et  $M_3$  = il est impossible :

$$a = M_1 [\text{âge} = [25,35]] \wedge M_2 [\text{diplôme} = \{A_1, A_2\}] \wedge M_3 [\text{nationalité} = \{\text{étranger}\}]$$

Un candidat se présente et répond aux caractéristiques suivantes : il a 27 ans, sous certaines conditions d'équivalence, il a un diplôme de type  $A_2$  et il n'est pas belge.

On le décrira par l'objet élémentaire où  $m_i$  exprime le degré de certitude :  $m_1$  = il est absolument vrai,  $m_2$  = il est vrai sous condition,  $m_3$  = il est absolument faux :

$$\omega^s = m_1 [\text{âge} = 27] \wedge m_2 [\text{diplôme} = \{A_2\}] \wedge m_3 [\text{nationalité} = \{\text{belge}\}]$$

La question est de savoir si  $\omega^s$  satisfait à l'assertion  $a$ . Pour le déterminer, on construit deux tables :

- 1 - la première permet dite de connexion permet de lier les  $M_i$  et les  $m_i$  pour obtenir les  $L_i$  exprimant un degré de convenance entre  $M_i$  et  $m_i$  et définir ainsi une fonction  $g$  dite de comparaison;
- 2 - la seconde,  $f$ , dite d'agrégation permet de faire la conjonction des  $L_i$ .

Grâce à la table de négation et au fait que les variables ne prennent qu'une seule valeur; par normalisation les événements élémentaires sont transformés pour qu'ils soient équivalents. Ainsi, en remplaçant dans  $\omega^s$ ,  $m_3 [\text{nationalité} = \{\text{belge}\}]$  par  $m_1 [\text{nationalité} = \{\text{étranger}\}]$  on se ramène à ce cas.

La table de connexion suivante permet de définir la fonction de comparaison  $g : m \times M \rightarrow L$  telle que  $g(m_i, M_j) = L_k$ . Dans l'exemple, on choisit :  $L = (L_1, L_2, L_3)$  avec  $L_1$  = convient,  $L_2$  = peut convenir et  $L_3$  = ne convient pas.



m	m <sub>1</sub>	m <sub>2</sub>	m <sub>3</sub>	m <sub>4</sub>
M				
M <sub>1</sub>	L <sub>1</sub>	L <sub>2</sub>	L <sub>2</sub>	L <sub>2</sub>
M <sub>2</sub>	L <sub>1</sub>	L <sub>2</sub>	L <sub>2</sub>	L <sub>2</sub>
M <sub>3</sub>	L <sub>3</sub>	L <sub>2</sub>	L <sub>1</sub>	L <sub>2</sub>
M <sub>4</sub>	L <sub>2</sub>	L <sub>2</sub>	L <sub>1</sub>	L <sub>2</sub>

Table de connexion :  $g(m_i, M_j)$

La table des agrégations est construite en posant  $f(L_i, L_j) = L_{\max(i,j)}$  qui convient à la sémantique de l'exemple où l'expert veut se montrer sévère.

L	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>
L			
L <sub>1</sub>	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>
L <sub>2</sub>	L <sub>2</sub>	L <sub>2</sub>	L <sub>3</sub>
L <sub>3</sub>	L <sub>3</sub>	L <sub>3</sub>	L <sub>3</sub>

Table des agrégations :  $f(L_i, L_j) = L_{\max(i,j)}$

Pour conclure notre exemple, nous avons :

$a(\omega) = g(m_1, M_1) \wedge g(m_2, M_2) \wedge g(m_1, M_3) = L_1 \wedge L_2 \wedge L_3 = L_3$  donc  $\omega$  ne convient pas, le candidat n'est pas retenu.

### 2.14.3. Les objets modaux de l'intérieur.

Les objets modaux de l'extérieur sont de la forme :  $a = \hat{I} M_i [y_i = V_i]$ , le mode  $M_i$  porte globalement sur l'événement  $[y_i = V_i]$  décrit par  $a$ . Les objets modaux de l'intérieur sont de la forme  $a = \hat{I} [y_i = M_i V_i]$  où le mode  $M_i$  concerne chaque  $V_i \subseteq O_i$ .

Pour définir un objet modal de l'intérieur, on se donne :

- $M^x$  : un ensemble de noms ou de nombres exprimant les modes associés à une connaissance du domaine de sémantique  $x$ . Exemple :  $M^x = [0,1]$ ;  $M^x = \{\text{rarement, parfois, souvent}\}$ .
- $Q_i = \{q_i^j\}_j$ , l'ensemble des applications  $q_i^j$  de  $O_i$  dans  $M^x$ .
- $y_i$  une variable, application de  $\Omega$  dans  $Q_i$ .
- $Op_x$  trois opérations ensemblistes définies dans  $Q_i$  :  $\cup_x, \cap_x, c_x$  qui sont respectivement des opérations d'union, d'intersection et de complémentation exprimant la sémantique  $x$ .
- $g_x$  est une application dite de comparaison de  $Q_i \times Q_i$  dans l'espace d'interprétation  $L_x$  ordonné et parfois identique à  $M^x$ .

- $f_x$  est une application symétrique, dite d'agrégation de  $P(L^x)$  (les parties de  $L^x$ ) dans  $L^x$ .

Définition :

Etant donnés  $Op_x$ ,  $g_x$  et  $f_x$ , un objet  $m_i$  est une application  $a_{YV}$  de  $\Omega$  dans  $L^x$  notée  $a = \hat{i}x$  [ $y_i = \{q_i^j\}_j$ ] telle que si  $\omega \in \Omega$  est décrit pour chaque  $i$  par  $y_i(\omega) = \{r_i^j\}_j \subseteq Q_i$  alors :

$$a_{YV}(\omega) = f_x(\{g_x(\bigcup_{j \in X} q_i^j, \bigcup_{j \in X} r_i^j)\}).$$

#### 2.14.4. Les objets modaux de l'intérieur associés à différentes sémantiques.

##### *Les objets possibilistes*

Les axiomes de base : on appelle « mesure de possibilité » une application notée  $\Pi$  de  $P(\Omega)$  dans  $[0,1]$  satisfaisant aux axiomes suivants :

- $\Pi(\Omega) = 1$ ;  $\Pi(\emptyset) = 0$ ;
- $\forall A, B \subseteq \Omega \quad \Pi(A \cup B) = \text{Max}(\Pi(A), \Pi(B))$ ;
- $\Pi(A) = 1 - N(\bar{A})$  où  $\bar{A} = c(A)$ , le complémentaire de  $A$  dans  $\Omega$ .

On peut appréhender intuitivement les notions par l'exemple suivant : soit  $\Pi_E(A)$  la possibilité pour que  $\omega \in \Omega$  soit dans  $A$  sachant qu'il est dans  $E$ . Cette possibilité est une application de  $P(\Omega)$  dans  $\{0,1\}$  elle vaut 1 si elle s'avère exacte et 0 sinon,  $\Pi_E(A)$  est nulle si  $A \cap E = \emptyset$  et vaut 1 si  $A \cap E \neq \emptyset$ . On a aussi  $N_E(A) = 1$  si  $E \subseteq A$  et  $N_E(A) = 0$  dans le cas contraire.

Les trois axiomes de base sont vérifiés.

La théorie des possibilités modélise différentes sortes de connaissances dont les trois suivantes :

- *la possibilité physique* : il s'agit d'exprimer la difficulté matérielle pour qu'une action puisse être effectuée, « un oiseau a la possibilité (physique) de voler »;
- *la possibilité au sens d'une concordance avec une connaissance actuelle* : surbase du bulletin météorologique, il est possible qu'il pleuve demain;
- *la possibilité exprimant le non-étonnement* : le mode de non-étonnement maximum a une probabilité égale à 1.

##### **Définition:**

Une assertion possibiliste est une assertion  $m_i$  qui prend ses valeurs dans  $L^P = [0,1]$ ; elle est définie par :

- $Op_p : \forall q_i^1, q_i^2 \in Q_i \quad q_i^1 \cup q_i^2 = \text{Max}(q_i^1, q_i^2); \quad q_i^1 \cap q_i^2 = \text{Min}(q_i^1, q_i^2)$ ;
- $\forall q \in Q_i \quad c_p(q) = 1 - q$ ;
- $g_p : g_p(q_i^1, q_i^2) = \text{Sup} \{ \min(q_i^1(v), q_i^2(v)) / v \in O_i \}$ ;
- $f_p : \forall L \subset L^P = [0,1]; \quad f(L) = \text{Max}\{l / l \in L\}$ .



On dit que les assertions qui viennent d'être définies sont possibilistes car elles se mettent sous la forme  $a = \hat{1}_p [y_i = \{q_i^j\}_j]$  et que les  $q_i^j$  sont des mesures de possibilité sur  $O_i$ .

### Les objets probabilistes

Les axiomes de base : étant donné  $C$  une algèbre d'événements sur  $\Omega$ , on appelle probabilité sur  $(\Omega, C)$  une application  $p$  de  $C$  dans  $[0,1]$  telle que  $P(\Omega) = 1$  est identifié à un événement dit certain d'extension  $\Omega$  et pour tout ensemble d'événements  $e_i$  d'extension  $a_i = |e_i|_\Omega$  disjointe, on a :

- $P(\bigcup_i a_i) = \sum_i P(a_i)$ .
- $P(\emptyset) = 0$ ,  $P(\bar{A}) = 1 - P(A)$ ,  $A \subseteq B$  implique  $P(A) \leq P(B)$ ,
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$ ,  $P(\bigcup_i a_i) \leq \sum P(a_i)$ .

La théorie des probabilités modélise différentes sortes de connaissances et au moins les suivantes :

*La chance* : à l'origine du calcul des probabilités, il s'agit de calculer la chance dans un jeu de hasard. La probabilité d'un événement est le rapport du nombre de cas favorables au nombre de cas possibles.

*La fréquence* : la probabilité représente une limite idéale de la fréquence d'un événement répété un très grand nombre de fois. La probabilité d'avoir pile est la limite du rapport : nombre de piles obtenus / nombre de tirages effectués.

*L'incertitude* : il s'agit de probabiliser des événements non répétables en mesurant un degré de certitude sur la réalisation d'un événement portant sur un objet unique : le prénom de la personne que je vais croiser est probablement François. L'incertitude provient d'un souvenir un peu effacé.

Définition :

Etant donné que chaque  $Q_i$  est un ensemble de mesures de probabilité sur  $O_i$ , nous avons la définition suivante :

Une assertion probabiliste est une assertion  $m_i$  qui prend ses valeurs dans  $L^{Pr} = [0,1]$ ; elle est définie par :

- $OP_{pr} : \forall q_i^1, q_i^2 \in Q_i \quad q_i^1 \cup_{pr} q_i^2 = q_i^1 + q_i^2 - q_i^1 q_i^2 ; q_i^1 \cap_{pr} q_i^2 = q_i^1 q_i^2 ;$
- $c(q_i^j) = 1 - q_i^j ;$
- $g_{pr} : \forall q_i^1, q_i^2 \in Q_i : g_{pr}(q_i^1, q_i^2) = \sum \{ \forall q_i^1(v) q_i^2(v) / v \in \Omega \} ;$
- $f_{pr} : f_{pr}(\{L_i\}) = \text{moyenne des } L_i.$

Exemple :

Un objet est décrit par sa couleur ( $y_1(w)$ ) (bleu ou rouge) et son format ( $y_2(w)$ ) (rond ou plat).

Soit  $a = [y_1 = q_i^1, q_i^2] \wedge_{pr} [y_2 = q_2]$  et  $\omega^s = [y_1 = \text{rouge}] \wedge_{pr} [y_2 = \text{rond}]$ .

Avec :

$q_1^1(\text{rouge}) = 0.9; q_1^1(\text{bleu}) = 0.1; q_1^2(\text{rouge}) = 0.5; q_1^2(\text{bleu}) = 0.5;$   
 $q_2(\text{rond}) = 0.2; q_2(\text{plat}) = 0.8.$

Ce qui signifie que a décrit des objets de deux types : soit souvent rouge et rarement bleu, soit rouge ou bleu avec une probabilité égale.

On obtient des objets dont la valeur est décrite par :

$$q_1^3(\text{rouge}) = 0.9 + 0.5 - 0.9 \times 0.5 = 0.95 \text{ et,}$$

$$q_1^3(\text{bleu}) = 0.1 + 0.5 - 0.1 \times 0.5 = 0.55.$$

D'autre part, on a :

$$r_1(\text{rouge}) = 1; r_1(\text{bleu}) = 0; r_2(\text{rond}) = 1; r_2(\text{plat}) = 0;$$

$$\begin{aligned} \text{On calcule : } a(\omega) &= g_{pr}(q_1^3, r_1) \wedge_{pr} g_{pr}(q_2, r_2) \\ &= (0.95 \times 1 + 0.55 \times 0) \wedge_{pr} (0.2 \times 1 + 0.8 \times 0) \\ &= 0.95 \wedge_{pr} 0.20 = f_{pr}(0.95 + 0.20) = 0.57. \end{aligned}$$

La probabilité pour que  $\omega$  satisfasse à a est 0.57.

### 3. Conclusion

La théorie des objets symboliques introduite par DIDAY permet d'étendre la problématique, les méthodes et les algorithmes de l'analyse classique des données à des données plus riches. Celles-ci expriment un niveau de connaissance plus élevé que celui fourni par de simples observations exprimables par un vecteur de valeurs quantitatives ou qualitatives.

Il s'agit de donner la possibilité d'utiliser en entrée, des données et des connaissances exprimées par des objets symboliques sans craindre de sortir du carcan tabulaire à n lignes (les individus) et p colonnes (les variables) et en évitant de perdre l'information par des modélisations ou codages arbitraires. On s'efforce d'exprimer les résultats sous forme d'objets symboliques, expressions obtenues automatiquement et possédant elles-mêmes un grand pouvoir explicatif.

Dans le cadre de ce premier chapitre, nous avons énoncé, défini et illustré les principales notions de la théorie des objets symboliques. Au sein du chapitre suivant, nous nous intéresserons à l'analyse symbolique de données symboliques.



## **1. Introduction**

L'aspect descriptif des données tel que nous venons de le présenter permet déjà d'exprimer beaucoup d'informations et d'en tirer des enseignements. Mais, au-delà de cette formalisation, nous voulons souligner les possibilités d'analyse que permettent les objets symboliques. Après avoir situé l'analyse des objets symboliques par rapport à l'analyse statistique classique et son intérêt dans différentes disciplines, nous détaillons deux études représentatives des potentialités de l'analyse symbolique.

En fin de chapitre, nous présentons un indice de dissimilarité permettant de comparer les objets entre eux.

## **2. L'approche symbolique : l'analyse des données**

Nous avons décrit un nouveau type de données, situons à présent l'analyse de ces données par rapport à l'analyse statistique classique.

### **2.1. Les quatre types d'analyse de données**

A partir de deux modes d'analyse et de deux sortes de donnée, DIDAY identifie quatre types d'analyse de données; entre elles, les frontières ne sont pas nécessairement marquées [DIDAY 93]. Nous présentons cette typologie avec le tableau 2.1.

	Données classiques	Données symboliques
Analyse Classique	(a)	(b)
Analyse Symbolique	(c)	(d)

Tableau 2.1. : typologie des analyses de données, tiré de [DIDAY 93].

- (a) → L'analyse numérique des données classiques : il s'agit ici du traitement des données quantitatives ou qualitatives avec des méthodes numériques fondées sur l'algèbre linéaire et utilisant des outils de la statistique (ACP, analyse discriminante, ...).
- (b) → L'analyse classique des données symboliques : pour réaliser des classifications ou une analyse factorielle, on utilise, par exemple, la distance euclidienne entre les objets.
- (c) → L'analyse symbolique des données classiques : dans ce cas-ci, on traite les tableaux de données classiques, c'est-à-dire des individus ou objets caractérisés par des variables quantitatives ou qualitatives, par l'approche symbolique. En utilisant les extensions, l'ordre symbolique, la généralisation, l'héritage, la qualité des objets, ... dès le départ sur les données classiques brutes ou après l'utilisation primaire d'analyse classique.
- (d) → L'analyse symbolique de données symboliques : on utilise, dans ce cas-ci, l'approche symbolique pour traiter des données qui sont aussi symboliques.



## **2.2. L'analyse des données symboliques par rapport à d'autres disciplines**

DIDAY dans [DIDAY 93] situe l'analyse des données symboliques, théorie dont il est l'initiateur, par rapport à d'autres disciplines. Il souligne ici l'intérêt de l'analyse symbolique pour ces disciplines.

En **Statistique** : élargissement du champ d'application à des populations qui sont généralement étudiées sous forme de points  $\mathbb{R}^p$  à des populations qui peuvent être formées d'objets complexes exprimant des sémantiques munies d'opérateurs non nécessairement numériques.

En **Intelligence Artificielle** : situé en amont des systèmes experts, il s'agit plutôt d'analyser des bases de connaissances ou de les induire à partir des données plutôt que d'étudier les inférences à partir de règles connues. En apprentissage, les problèmes traités sont proches de ceux de l'analyse des données, mais s'en distinguent par les objets traités et les méthodes.

Par rapport à la **logique floue**, les objectifs se situent à des niveaux différents puisque le but de l'analyse symbolique est de faire de l'analyse de données. On peut néanmoins utiliser des fonctions floues pour définir des objets symboliques. C'est le cas des objets possibilistes : d'autres axiomes que ceux du flou sont utilisés dans le cas des objets probabilistes ou crédibilistes.

## **2.3. Les six étapes d'une analyse de données symboliques [DIDAY 93]**

En général, lorsque l'on travaille avec un ensemble de nombreux individus, les objectifs des analystes sont plutôt descriptifs, DIDAY fournit une séquence logique des opérations à réaliser.

- (1) ↪ Partir d'un ensemble d'objets complexes,
- (2) ↪ extraire des classes par classification, analyse factorielle, arbres de décisions, treillis, ...
- (3) ↪ représenter ces classes afin d'obtenir des objets définis en intentions,
- (4) ↪ construire des objets symboliques permettant d'identifier des objets individuels,
- (5) ↪ analyser, synthétiser, classier, discriminer, organiser par différentes méthodes de l'analyse des données symboliques l'ensemble des objets symboliques issus de l'étape (4),
- (6) ↪ déduire de ces analyses des métaconnaissances telles que une pyramide d'héritage et des règles entre les objets symboliques.



### 3. L'approche symbolique : exemples d'application

Les analyses de données symboliques s'adressent à des données très variées. Nous avons choisi de détailler deux études : l'une concernant la sécurité automobile, l'autre concernant les tactiques de pêche. Ces deux analyses sont très différentes, elles expriment le potentiel de l'approche symbolique. Dans le cadre du mémoire, nous ne tenons pas à montrer toutes les possibilités de ce type d'analyse, mais à en illustrer quelques aspects.

#### 3.1. Scénarii d'accidents : un outil pour les diagnostics de sécurité

Nous étudions ici l'analyse réalisée par REGNIER en 1992 [REGNIER 92]. Afin d'améliorer la sécurité routière, les spécialistes disposent de grandes banques de données constituées des procès verbaux d'accident rédigés par les gendarmes de terrain. Pour éviter la reproduction des accidents, des experts (FLEURY, FLIN, PEYTAVIN) de l'INRETS (Institut de Recherche sur les Transports et leur Sécurité), ont eu l'idée en 1991, de formuler des scénarii d'accidents afin d'aider à la mise en oeuvre de mesures appropriées. Les scénarii des experts, définis sous forme de phrases décrivent des caractéristiques d'usagers, de déplacements, de lieux, de périodes et parfois de véhicules.

Exemple de scénario d'accident :

« homme de 30-50 ans perdant le contrôle de son véhicule (usager local, expérimenté, souvent sous influence de l'alcool), accident survenant le jour »

Remarquons que jusqu'ici, les objets symboliques n'interviennent pas: les experts ont simplement exprimé de façon plus claire et simple les procès verbaux. Mais l'expression se rapproche tout de même (intuitivement) de la formalisation sous forme d'assertion.

C'est REGNIER dans [REGNIER 92] qui, dans le cadre de son rapport de stage, utilise l'ensemble de ces données comme base pour son analyse.

L'auteur a appliqué plusieurs étapes de l'analyse des données symboliques pour améliorer, compléter et organiser la base de scénarii d'accidents fournie par les experts. Dans ce travail, chaque accident est considéré comme « intention » d'une classe d'accidents de la base de données de la Gendarmerie. Les étapes de l'étude ont été les suivantes :

a) Exprimer les scénarii sous forme d'objets symboliques, en effet, ce sont les objets symboliques qui se sont avérés les plus adéquats.

Exemple :

scénario = [jour = {70% lundi, 30% dimanche}]  
           $\wedge$  [état de la surface = {enneigée, verglassée}]  
           $\wedge$  [signalisation = {60% (stop), 40% (cédez le passage, priorité à droite)}]



b) Calculer leur extension dans la base des prototypes.

Exemple :

Si un scénario est décrit par

$\text{scen} = [\text{jour} = \{70\% \text{ lundi}, 30\% \text{ dimanche}\}] \wedge [\text{heure} = [7,9]]$ ,

et si un accident de la base est défini par

$\text{acc} = [\text{jour} = \{\text{lundi}\}] \wedge [\text{heure} = 8] \wedge [\text{mois} = \text{février}]$ ;

alors, d'après la définition des objets probabilistes, le degré d'appartenance de cet accident à ce scénario est calculé par

$\text{scen}(\text{acc}) = (70 \times 1 + 30 \times 0 + 100 \times 1 + 100 \times 1) / 3 = 90$ .

Justification :

(-)  $70 \times 1 + 30 \times 0$  car le jour est un lundi et pas un dimanche,

(-)  $100 \times 1$  car l'accident se déroule bien dans l'intervalle horaire du scénario,

(-)  $100 \times 1$  car [mois = quelconque] n'est pas indiqué dans le scénario, mais reste sous-entendu.

On peut calculer une extension à un seuil donné en deçà duquel les accidents sont rejetés. Le seuil peut être fourni par l'expert au vu des résultats et de son savoir, ou bien calculé automatiquement en faisant la somme des plus petites probabilités associées à chaque événement. De la même façon, on peut associer un seuil maximum en prenant les probabilités les plus importantes. Ce double calcul des seuils permet de définir des prototypes au sens de la terminologie utilisée en sciences cognitives.

Exemple :

Sur base des données citées plus haut, on peut calculer :

- un seuil minimum :  $(30 + 100) / 2 = 65$  et,

- un seuil maximum :  $(70 + 100) / 2 = 85$ .

$\text{proto} = [\text{jour} = \text{lundi}] \wedge [\text{heure} = 8]$  est un prototype.

c) Amélioration des scénarii fournis par les experts.

Certains scénarii s'avèrent d'extension trop grande ou trop petite. Dans le premier cas l'expert doit ajouter des événements spécifiques, dans l'autre cas il doit en supprimer. Il peut être aidé dans cette tâche, par le calcul de l'extension de ces événements.

d) Construction de nouveaux scénarii.

Les experts disposaient de 579 accidents. A partir de ces informations, 12 scénarii ont été finalement constitués. L'union de leurs extensions (des scénarii exprimés sous forme de phrases) couvrait 286 accidents. Il restait donc 286 accidents non couverts, auxquels il fallait associer de nouveaux scénarii. Une classification automatique de ces accidents fournit 12 classes dont les intentions ont été proposées aux experts. Après discussions, 26 scénarii ont été retenus. Leurs extensions couvrent 93,4% des accidents du fichier de départ.

e) Organisation des scénarii.

Les étapes a), b), c), d) permettent de construire une base de connaissance définie par un ensemble de 26 objets symboliques. Ces étapes correspondent aux quatre premières phases décrites au point 2.3. ci-avant. L'étape suivante consiste à organiser les objets obtenus (étape (5) du point 2.3.).

Une pyramide d'héritage a été construite de façon ascendante. Nous donnons un bref aperçu sur la figure 2.1. correspondant aux accidents de collision entre deux usagers.



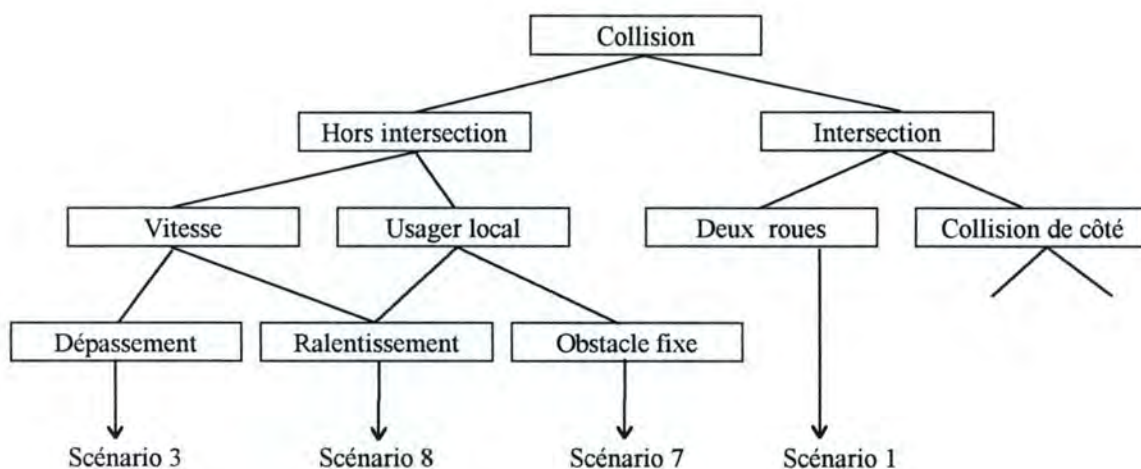


Figure 2.1. : Organisation des scénarii par une pyramide d'héritage, tiré de [REGNIER 92].

#### f) Connaissances sur les connaissances.

Cette étape finale consiste à interpréter les objets symboliques et la pyramide d'héritage. Les objets symboliques sont partitionnés en deux groupes dits forts ou faibles. Sont considérés comme forts ceux qui, malgré une description longue (comprenant beaucoup d'événements élémentaires [ $y_i = q_i$ ]), ont une grande extension, les autres sont qualifiés de faibles. L'auteur s'est aperçu que les scénarii établis par les experts de l'INRETS sont en majorité forts, ce qui valide leur étude.

Dans ce cas-ci, l'analyse des données symboliques a permis de représenter, confirmer, compléter et organiser les scénarii d'accidents issus de la base de données. On peut mieux différencier les types d'accidents et spécifier les connaissances apportées par les scénarii.

### 3.2. Analyse symbolique des tactiques de pêche artisanale au Sénégal

Avec plus de 200.000 tonnes de poissons débarqués chaque année, le secteur de la pêche artisanale au Sénégal joue un rôle de premier plan dans l'économie du pays. Caractérisé par sa complexité et sa grande adaptabilité face aux instabilités environnementales et socio-économiques, le système de pêche artisanale demande à être mieux cerné, mieux compris dans l'optique d'une gestion plus efficace de la richesse halieutique. C'est dans ce cadre que PERINEL situe son étude en 1992 [PERINEL 92]. L'auteur a choisi de s'intéresser à un aspect essentiel de la pêche artisanale : le comportement tactique et stratégique du pêcheur. Le but étant de mieux cerner les décisions et réactions des pêcheurs face aux fluctuations de leur environnement.

L'approche globale est schématisée dans la figure 2.2. suivante :

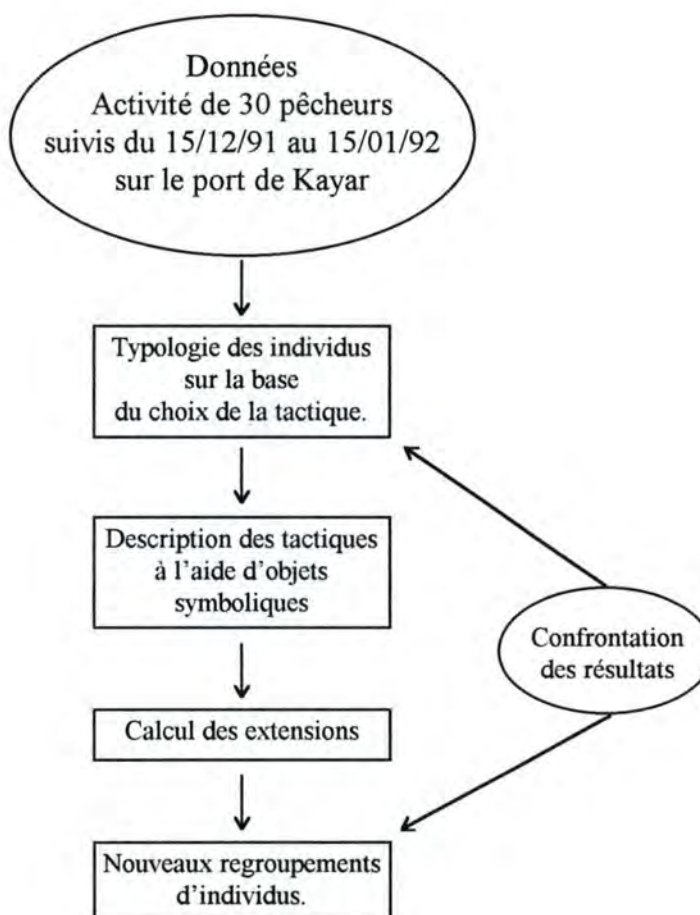


Figure 2.2. : Schéma d'analyse, extrait de [PERINEL 92].

Dans le cadre de l'approche numérique du problème, l'auteur désire obtenir des groupes homogènes d'individus en regard de la tactique de pêche adoptée. Les variables retenues, relatives aux différentes décisions du pêcheur avant sa sortie en mer, sont les suivantes :

- le type d'engin,
- l'effectif de l'équipage,
- les espèces ciblées,
- le lieu de pêche.

L'utilisation de méthodes d'analyse de données multidimensionnelles et en particulier de la classification retient l'attention de l'auteur. Deux étapes ont permis de réaliser la typologie recherchée :

- 1 - détermination d'une typologie basée sur les espèces de poissons ciblées et,
- 2 - détermination des tactiques de pêche proprement dites.

Résultats de la première étape :

La pêche artisanale sénégalaise est caractérisée par un aspect multispécifique (certains pêcheurs ciblent plusieurs espèces lors d'une même sortie). Pour ne pas que la classification soit trop influencée par les quatre variables représentant les espèces ciblées, l'auteur procède à



une première classification permettant la création d'une nouvelle variable synthétique résumant l'appartenance de chaque individu à une classe d'espèces ciblées.

C'est par l'utilisation de variables binaires indiquant pour chaque individu la présence ou l'absence de l'espèce comme cible que l'auteur a résolu le problème. La typologie des individus sur la base des espèces ciblées (17 ≠) est réalisée via une analyse des correspondances multiples suivie d'une classification hiérarchique ascendante. Le dendrogramme, résultat graphique de la classification, incite PERINEL à retenir 8 classes.

Résultats de la seconde étape :

A l'aide des résultats de la première étape, de l'ajout des variables portant sur le lieu de pêche, l'équipage embarqué, le type d'engin, PERINEL détermine les tactiques de pêche proprement dites. Méthodologiquement, la procédure est identique à celle de la première étape : une analyse des correspondances multiples puis une classification hiérarchique ascendante. Les résultats permettent d'identifier cinq types de tactiques de pêche.

Concernant l'approche symbolique du même problème, la première étape est de formuler sous forme d'assertions probabilistes les différentes tactiques identifiées lors de la phase précédente. En effet, grâce à l'analyse numérique, l'auteur a pu dégager au sein du port de Kayar, et ce pendant une période donnée, cinq groupes homogènes d'individus, chacun de ces groupes reflétant un comportement tactique type.

Exemple de tactique formalisée :

tactique 1 = [espèce cible 1 = 125]  
 ^ [espèce cible 2 = aucune]  
 ^ [engin = {70% (type 2), 30% (type 3)}]  
 ^ [lieu = {28% (13), 18% (14), 11% (21), 17% (7), 26% (autre)}]

Lors du calcul des extensions, deux indices seront calculés :

- *un indicateur du degré de généralisation.*

On cherche à retrouver dans l'extension d'un objet tactique  $t_i$ , le maximum d'individus appartenant à la classe  $C_i$  dégagée de l'analyse numérique.

$$G_i = \frac{\text{card}\{ |t_i|_{C_i} \}}{\text{card}\{C_i\}}$$

$$= \frac{\text{le nombre d'individus dans l'extension de } t_i \text{ appartenant à la classe } C_i}{\text{le nombre d'individus appartenant à la classe } C_i}$$



**- un indicateur exprimant le pouvoir de discrimination.**

Une tactique  $t_i$  doit discriminer au mieux les individus des différentes classes  $C_i$ .

$$D_i = \frac{\text{card}\{ |t_i|_{C_i} \}}{\sum_j \text{card}\{ |t_i|_{C_j} \}}$$
$$= \frac{\text{le nombre d'individus classés } C_i \text{ appartenant à l'extension de } t_i}{\text{le nombre total d'individus de l'extension}}$$

Le calcul des extensions conduit à la formation de regroupements d'individus. A la différence des résultats d'une classification (présentée ci-avant), il ne s'agit pas d'une partition des objets observés puisqu'un même individu peut se trouver affecté à plusieurs extensions simultanément. En plus de cette analyse, l'auteur a repris la phase de calcul des extensions pour l'appliquer cette fois à un groupe de pêcheurs n'ayant pas participé à l'analyse. Certains ne sont pas affectés à une extension (11/200) tandis que d'autres appartiennent à deux extensions (19/200).

Ayant identifié différentes tactiques, PERINEL étudie alors le comportement stratégique des pêcheurs. Il reprend dans un tableau, et ce pour 13 pêcheurs, la tactique choisie chaque jour pendant un mois. Chaque ligne du tableau décrit pour un pêcheur son profil tactique, c'est-à-dire sa stratégie, sa façon de coordonner ses actions. L'auteur réalise une analyse factorielle et observe le regroupement de profils similaires.

Au terme de son étude, PERINEL conclut que l'analyse des objets symboliques vient combler certaines lacunes du numérique en apportant du sens, de l'explicatif, en utilisant les connaissances du domaine. L'apport s'est révélé surtout intéressant dans le cadre de son étude pour formaliser de nouveaux concepts que sont les tactiques de pêche exprimées sous formes d'objets explicites. Le but de l'analyse est aussi de tenter de mieux comprendre les raisons qui amènent les pêcheurs à changer de tactique et par conséquent de mettre en évidence d'éventuelles règles de changement tactique.

#### **4. Dissimilarité.**

Nous venons de décrire des analyses qui utilisent plutôt les classifications, les regroupements. Lorsque l'on se place à un niveau plus microscopique, si l'on s'intéresse à la comparaison d'objets deux par deux, une série d'indices ont été développés pour permettre des comparaisons. Nous verrons dans le chapitre suivant une méthode de comparaison visuelle : la représentation graphique.

Nous reprenons ici les conclusions des travaux de DE CARVALHO dans ce domaine. Dans le cadre de sa thèse, le chercheur a étudié les indices qui permettent de différencier les objets symboliques entre eux. Après plusieurs améliorations successives et l'étude de la littérature dans le domaine, l'auteur présente un indice de dissimilarité entre objets symboliques.



Francisco DE CARVALHO dans [DE CARVALHO 94] explique comment un indice peut exprimer une dissimilarité entre deux objets symboliques. Sur bases de différentes études publiées et de ses recherches sur d'autres indices, il présente un indice de dissimilarité entre objets symboliques booléens basé sur leurs extensions. Dans le cadre du mémoire, nous reprenons ses principaux résultats, le lecteur intéressé trouvera dans [DE CARVALHO 94] le détail des propriétés et leurs démonstrations.

Dans les analyses statistiques classiques, lorsque les individus sont numériques, des indices basés sur la moyenne, la médiane, ... peuvent calculer une forme de similarité. Mais, nous l'avons vu dans le chapitre précédent, un objet symbolique est décrit par plusieurs variables, une variable peut être quantitative (discrète ou continue) ou qualitative (ordinaire ou nominale) et elle peut ne prendre aucune valeur ou en prendre une ou plusieurs pour un Objet Assertion Booléen (*OAB*). Pour représenter des connaissances réelles, la description d'une classe d'individus par un *OAB* doit tenir compte de différents types de dépendances logiques (*DL*) entre variables. Par exemple, nous ne pouvons pas décrire la couleur du chapeau d'une espèce de champignon qui n'a pas de chapeau (dépendance conditionnelle, *DC*). Et un expert peut savoir que si la couleur du chapeau d'une espèce de champignon est blanche alors sa taille est inférieure à 5 cm (dépendance de corrélation logique, *DCL*). Ces *DL* sont exprimées par des règles entre les variables.

Lors de la description d'une classe d'individus par les objets usuels de l'analyse de données, la variabilité est prise en compte par la moyenne, le mode, l'écart type, etc., et non par une disjonction de valeur relative à une variable. Certains auteurs, tels que PEARSON, MAHALANOBIS, SANGHVI et CROVELLO dans [DE CARVALHO 94], ont proposé des indices qui tiennent compte de la variabilité exprimée de cette façon lors du calcul de la proximité entre ces objets. D'autres, tels KENDRICK et GOWER dans [DE CARVALHO 94], ont proposé des approches pour tenir compte de l'influence de la *DC* entre variables lors du calcul de la proximité entre les objets. Enfin, JARDINE et SIBSON dans [DE CARVALHO 94] ont proposé une approche (*J - Dissimilarité*) pour tenir compte en même temps de l'influence de la variabilité (exprimée par la moyenne, etc.) et de la *DC* entre variables lors du calcul de la proximité entre objets.

#### **4.1. Le potentiel de description d'un objet assertion booléen.**

Soit  $a_j = [y'_1 = V_1^j] \wedge \dots \wedge [y'_q = V_q^j]$  un *OAB*, où  $e_i^j = [y'_i = V_i^j]$  est le  $i^{\text{ème}}$  événement élémentaire booléen (*EEB*). L'auteur définit le *Potentiel de Description (PD)* d'un *OAB*  $a_j$  noté  $\pi(a_j)$  comme :

$$\pi(a_j) = \prod_{i=1}^q \pi(e_i^j)$$



Où  $\pi(e_i^j)$  est le *PD* de l'*EEB*  $e_i^j$ , qui est :

$$\pi(e_i^j) = \mu(V_i^j) = \begin{cases} \text{cardinal}(V_i^j), & \text{si } y_i \text{ est discrète ou,} \\ \text{écart}(V_i^j), & \text{si } y_i \text{ est continue.} \end{cases}$$

Où :

- $\text{cardinal}(V_i^j)$ , représente le nombre d'occurrence que prend la variable,
- $\text{écart}(V_i^j)$ , la différence entre les limites supérieure et inférieure d'un intervalle.

Cette formule est applicable facilement lorsque il n'y a pas de *DL* entre les *OAB*. L'auteur propose un élargissement de la méthode pour le calcul du *PD* pour des *OAB* pour lesquels il existe une *DL*.

Les *DL* entre variables peuvent être représentées par un ensemble de graphes connexes, où les noeuds représentent les variables. L'axe fléché détermine la dépendance (dans l'exemple de la figure 2.3. :  $y_1$  influence les variables  $y_2$  et  $y_3$ ). L'étiquette exprime le type de dépendance (1 si la règle exprime une *DC* et 2 si la règle exprime une *DCL*).

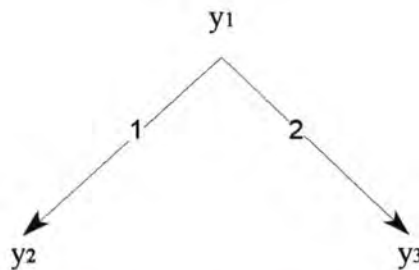


Figure 2.3. : Graphe de *DL* entre variables, extrait de [DE CARVALHO 94] modifié.

Pour permettre le calcul de *PD*, on individualise tous les graphe connexes, la formule est alors appliquée sur les *OAB*, la procédure est répétée sur chaque *OAB*.

#### 4.2. Le calcul de la proximité entre objets assertion booléens.

Soit  $a_k$  et  $a_j$  deux *OAB*, dans cette approche, la proximité de deux *OAB* est basée sur leur *PD*. Plus précisément, l'idée centrale de cette approche, inspirée de la figure 2.4. ci-dessous, est l'hypothèse que la dissimilarité entre deux *OAB* est fonction du nombre de descriptions d'individus propres à chaque *OAB* ( $b$  et  $c$ ) et du nombre de descriptions qui ne sont de ni l'un ni l'autre *OAB* ( $d$ ).



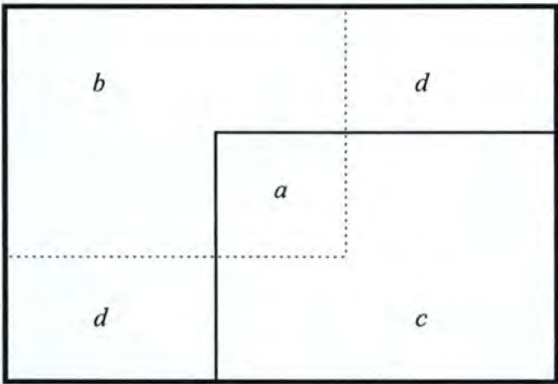


Figure 2.4. : représentation cartésienne de deux *OAB* (les deux petits rectangles intérieurs) et de leur union symbolique, extrait de [DE CARVALHO 94] et modifié.

On peut exprimer cette représentation aussi sous la forme suivante :

$a_j \backslash a_k$	Accord	Désaccord	Total
Accord	$A = \pi(a_j \cap_a a_k)$	$B = \pi(a_j) - \pi(a_j \cap_a a_k)$	$A + B = \pi(a_j)$
Désaccord	$C = \pi(a_k) - \pi(a_j \cap_a a_k)$	$D = N - A - B - C$	$C + D = \pi(a_j \cup_a a_k) - \pi(a_j)$
Total	$A + C = \pi(a_k)$	$B + D = \pi(a_j \cup_a a_k) - \pi(a_k)$	$N = A + B + C + D = \pi(a_j \cup_a a_k)$

Tableau 2.2. : comparaison entre  $a_j$  et  $a_k$ , extrait de [DE CARVALHO 94].

Où :

- A mesure le nombre de descriptions d’individus partagés par  $a_j$  et  $a_k$ (région a);
- B mesure le nombre de descriptions d’individus propres à  $a_j$  (région b);
- C mesure le nombre de descriptions d’individus propres à  $a_k$  (région c);
- D mesure le nombre de descriptions qui ne sont dans aucune de celles des *OAB*.

Compte tenu de la figure 2.4. et du tableau 2.2. l’auteur choisit comme indice de proximité:

$$d_{ext}(a_j, a_k) = B + C + D = \pi(a_j \cup_a a_k) - \pi(a_j \cap_a a_k)$$

Dans le cadre de ses publications, DE CARVALHO énonce et démontre une série de propositions et de propriétés. Nous n’en parlons pas dans ce travail, elles ne se rapportent pas directement au sujet.

Illustrons par un exemple simple :

$$a = [y_1 = [10,50]] \wedge [y_2 = [100,400]] \wedge [\text{Couleur} = \{\text{blanc, rouge}\}]$$

$$b = [y_1 = [30,50]] \wedge [y_2 = [200,500]] \wedge [\text{Couleur} = \{\text{rouge}\}]$$

$$c = [y_1 = [5,30]] \wedge [y_2 = [50,600]] \wedge [\text{Couleur} = \{\text{vert, noir, rouge}\}]$$

Calculons l'indice pour les objets a et b :

Etape	Développement	Calcul	Résultat
$\bigcup_{a,b}$	$[10, 50] \wedge [100, 500] \wedge [\{\text{blanc, rouge}\}] =$	$40 \times 400 \times 2 =$	32000
$\bigcap_{a,b}$	$[30, 50] \wedge [200, 400] \wedge [\{\text{rouge}\}] =$	$20 \times 200 \times 1 =$	4000
$d_{\text{ext}}(a,b) =$	$\bigcup_{a,b} - \bigcap_{a,b} =$	$32000 - 4000 =$	28000

Pour les objets a et c :

Etape	Développement	Calcul	Résultat
$\bigcup_{a,c}$	$[5, 50] \wedge [50, 500] \wedge [\{\text{blanc, rouge, vert, noir}\}] =$	$45 \times 550 \times 4 =$	99000
$\bigcap_{a,c}$	$[10, 30] \wedge [100, 400] \wedge [\{\text{rouge}\}] =$	$20 \times 300 \times 1 =$	6000
$d_{\text{ext}}(a,c) =$	$\bigcup_{a,b} - \bigcap_{a,b} =$	$99000 - 6000 =$	93000

L'interprétation :

Puisque  $d_{\text{ext}}(a,c)$  est plus élevé que  $d_{\text{ext}}(a,b)$ , on peut dire que c est plus différent de a que ne l'est l'objet b. Plus la valeur de  $d_{\text{ext}}$  est élevée, plus les objets sont « dissimilaires », différents selon l'indice. Deux objets identiques ont un  $d_{\text{ext}}$  nul.

Remarque :

Dans les applications réelles, où le nombre de variables est important, pour éviter de manipuler des nombres trop grands, DE CARVALHO propose d'utiliser le  $\ln$  de  $d_{\text{ext}}$  lorsque celui-ci est  $> 0$ .



## 5. Conclusions [DIDAY 93].

L'analyse des données symboliques trouve sa place dans un créneau entre l'intelligence artificielle et la statistique, entre les approches logiques, symboliques et numériques. Elle ouvre la voie à un grand champ d'applications : celui du traitement des objets complexes en tenant compte de connaissances non nécessairement d'ordre purement numérique.

Les applications ayant donné des résultats intéressants touchent des domaines aussi divers que les signaux radars, les stratégies de pêche, des scénarii d'accident, des maladies ou des espèces de fleurs. Dans le cadre de notre travail, nous nous sommes limités à deux illustrations.

Les deux exemples que nous avons développés permettent de lever un coin du voile sur les potentialités offertes par l'analyse symbolique des données symboliques vis à vis de l'analyse classique de données numériques. L'analyse symbolique comble un créneau situé entre l'Intelligence Artificielle et la Statistique.

Par l'analyse de scénarii d'accidents, l'analyse symbolique ouvre la voie au traitement d'objets complexes en tenant compte de connaissances non nécessairement d'ordre purement numérique. L'analyse des données symboliques a permis de représenter, confirmer, compléter et organiser les scénarii d'accidents issus d'une base de données; les chercheurs ont pu différencier les types d'accidents et spécifier les connaissances apportées par les scénarii.

La seconde analyse résumée rapproche l'analyse des données de l'Intelligence Artificielle. La méthodologie présentée par PERINEL permet de montrer que les deux approches (numérique et symbolique) sont intégrées de façon conjointe et complémentaire :

- L'approche numérique est caractérisée par son efficacité, les méthodes de classification sont utilisées dans la première phase pour déterminer une partition en classes homogènes des pêcheurs.

- L'approche symbolique vient combler certaines lacunes du numérique pur en apportant du sens, de l'explicatif en utilisant les connaissances d'un domaine. Dans le cadre de l'étude, l'apport s'est révélé intéressant pour formaliser de nouveaux concepts que sont les tactiques de pêche exprimées sous forme d'objets explicites (utilisation de la logique modale).



### 1. Introduction

Nous venons de le voir : l'objet symbolique est un concept qui permet de décrire de façon assez complète des individus complexes. Que ce soit sous forme de tableaux ou d'assertions, percevoir l'information véhiculée n'est pas nécessairement chose aisée : une représentation graphique s'impose.

La complexité de l'objet symbolique est telle qu'il est nécessaire d'opter pour une représentation propre. L'idéal est de représenter un maximum d'informations sans surcharges. Sur le marché, il n'existe pas de logiciels permettant ce type de représentation (multivariée, multi-échelle, ...)<sup>1</sup>. Dans le cadre de leurs recherches, M. NOIRHOMME-FRAITURE et M. ROUARD ont développé une solution pour la représentation d'objets statistiques complexes : « l'étoile zoom » [NOIRHOMME-FRAITURE 96], [NOIRHOMME-FRAITURE 97a] et [NOIRHOMME-FRAITURE 97b].

### 2. L'Etoile Zoom

Le principe de la représentation en étoile, c'est-à-dire, où les axes (souvent plus de 4) ont même origine et sont répartis régulièrement sur un éventail de 360°, est déjà utilisé dans d'autres disciplines. Par exemple en climatologie, la rose de vents schématise, pour une période donnée, les fréquences des vents en fonction de leurs directions.

Chaque axe correspond, ici, à une direction cardinale (le Nord toujours au-dessus). Nous présentons à la figure 3.1 ci-contre, un exemple dont les données sont tirées de « Géographie de la Belgique » [DENIS 92] : pourcentages moyens des vitesses de vents mesurées à 1,5 m. du sol à la station météorologique de Uccle de 1901 à 1930.

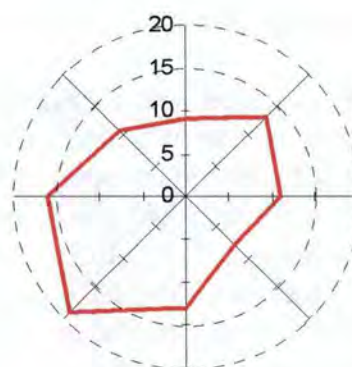


Figure 3.1 : Rose des vents.

En géologie, l'utilisation de graphiques radiaires permet de classer les différentes roches selon leurs compositions minéralogiques. Un axe représente le pourcentage de présence d'un composant minéralogique ( $\text{CaCO}_3$ ,  $\text{SiO}_2$ ,  $\text{NaCl}$ , ...). Chaque roche a une composition type, sa « signature ». Lors de l'analyse d'un nouvel échantillon, on dessine sa signature et on compare aux signatures connues des différentes roches; par comparaison visuelle, on détermine le type de roche auquel l'échantillon appartient.

<sup>1</sup> Vérifications faites auprès des logiciels statistiques « courants » : Lotus 5.0, Excel 5.0, STATISTICA, Glim 4, SAS 6.11, SPSS, Intercooled Stata 4.0, ...



Un troisième exemple d'utilisation, en géomorphologie, citons les roses d'orientation des cailloux dans les dépôts détritiques. L'orientation d'un élément (un cailloux) est la valeur de l'angle que fait le plus grand axe du cailloux avec la direction Nord. Les résultats des mesures sont consignés sur des axes d'orientation groupant les valeurs par secteurs de 10°. L'examen de ces roses d'orientation détermine le sens d'écoulement d'une paléo-rivière. PISSART souligne dans [PISSART 92] les intérêts d'une telle représentation : la rapidité de réalisation et l'apport visuel d'un tel graphique. Ce travail ne nécessite pas l'emploi d'outils perfectionnés et peut être réalisé sur le terrain.

Dans le cas de l'exemple de la figure 3.1. ci-avant, tous les axes ont même échelle et représentent une même variable pour une direction donnée. Dans le cas des graphiques radiaires de classification géologique, chaque axe correspond à un pourcentage de présence minéralogique; dans le cas des objets symboliques, l'échelle de représentation d'un axe à l'autre peut varier. On peut imaginer une représentation en étoile où chaque axe porte une variable différente, selon une échelle différente. Citons par exemple les diagrammes de KIEVIET.

## **2.1. Introduction à l'Etoile Zoom**

Cette présentation se base sur les études publiées par M. NOIRHOMME-FRAITURE et M. ROUARD [NOIRHOMME-FRAITURE 96], [NOIRHOMME-FRAITURE 97a] et [NOIRHOMME-FRAITURE 97b].

Les hypothèses de travail sont les suivantes :

- classiquement, en analyse de données, les individus sont caractérisés par les lignes d'un tableau où les colonnes sont les variables;
- en statistique, pour imager de grandes quantités d'informations, l'utilisation de représentations graphiques est largement répandue;
- par nature, les objets symboliques sont complexes à synthétiser, à représenter;
- les types de variables, au sein d'un objet symbolique sont très différents (quantitatif, qualitatif, intervalles, ordinal ou nominal, ...);
- le nombre de variables peut varier d'une situation à une autre;
- les domaines d'étude sont très variables : scénario d'accidents de la circulation, tactiques de pêche, ...

De ce constat, les auteurs soulignent la nécessité de concevoir une méthode de représentation simple, sans être simpliste, permettant de visualiser toute l'information nécessaire sans surcharger le graphique.

Un prototype de visualisation est développé, il sera intégré à un logiciel plus complet. A terme, ce logiciel développé dans le cadre d'un projet européen, est appelé à gérer la représentation d'objets symboliques. La représentation graphique est nécessairement synthétique. C'est une image synthétique, et donc partielle, de l'objet. Le graphique et le tableau de nombres sont complémentaires. Le graphique fournit un cliché, une synthèse aisément interprétable du tableau complet de nombres et souvent trop peu lisible. Le graphique permet donc une meilleure compréhension et souvent une meilleure comparaison visuelle entre individus.



2.2. La représentation des objets symboliques

Les objets symboliques se caractérisent par une complexité de représentation; celle-ci est liée à la richesse des informations que ces objets véhiculent par rapport à des données statistiques classiques.

Les objets symboliques, appelés également assertions, sont décrits par un ensemble de propriétés ou variables qui pourront être de types différents. Si, par exemple, l'on s'intéresse aux objets « voiture d'occasion », les types de variables pourraient être les suivants :

- *variables quantitatives* (valeur entière ou réelle), exemple : Cylindrée = {1900 cm<sup>3</sup>};
- *variables qualitatives*, (ordinales, exemple : Mois de mise en circulation = {janvier, février, mars, avril, mai juin, juillet, août, septembre, octobre, novembre, décembre} ou nominale, exemple : carburant = {diesel, essence});
- *variables en intervalle*, exemple : Prix de vente (kBEF) = [150; 160] signifie que le prix de vente est compris entre 150.000 et 160.000 BEF;
- *valeurs pondérées*, par exemple : Etat = [accidenté (0,25), intact (0,75)] signifie que 25% de la population a « accidenté » comme état et 75% sont des véhicules intacts.
- *variables multivaluées*, par exemple : Catégorie {citadine (C), familiale (F), monovolume (M), sportive (S)} ce qui signifie qu'un élément de la population peut être soit C, soit F, soit M, soit S. Soit éventuellement plusieurs, cela est permis.

Nous reprenons à la figure 3.2., un exemple présentant les types décrits ci-dessus.

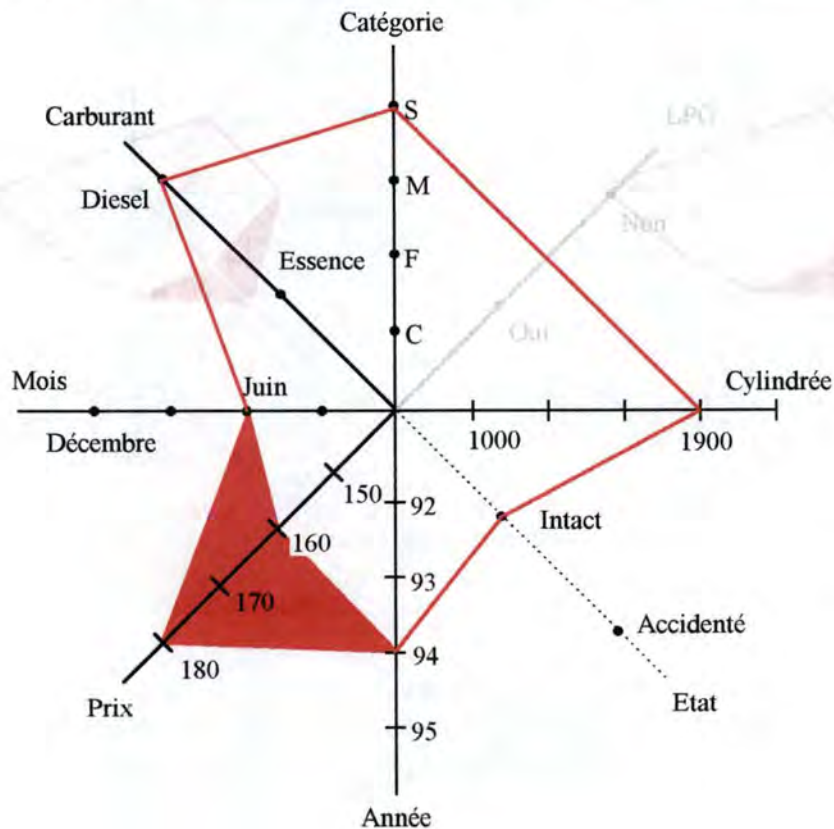


Figure 3.2. : exemple de représentation d'un objet symbolique avec la méthode de l'étoile zoom.



L'avantage majeur de cette représentation en étoile est tout d'abord de délivrer une image synthétique de l'objet. Chaque axe correspond à une variable, ce qui permet de représenter dans ce cas, huit variables différentes.

- Une variable quantitative est représentée sur un axe gradué. (exemple : cylindrée moteur ( $\text{cm}^3$ ), l'année de fabrication ).
- Les catégories d'une variable qualitative sont représentées par des points sur un axe, points distribués équitablement, harmonieusement sur l'axe. (exemples : état du véhicule, le pays d'origine (Espagne (E), Japon (J), Allemagne (D), France (F) ).
- Les valeurs de chaque variable sont reliées afin de confectionner une étoile. Si un intervalle est utilisé, ses limites sont liées et toute la surface est colorée (exemple : prix d'achat).
- Lorsqu'un objet n'a pas de valeur pour une variable, l'axe est dessiné en gris clair et la forme de l'étoile résultante ne tient pas compte de cet axe. (exemple : présence de la propulsion au LPG).
- Les probabilités n'apparaissent pas directement sur le graphique 2D dans le but de simplifier au maximum la représentation, par conséquent, seule la valeur de plus grand poids est représentée sur l'étoile. Ces axes sont représentés en pointillés. En cliquant sur l'axe, un graphique 3D des pondérations est affiché (exemple : état du véhicule). Ce cas est illustré à la figure 3.3. ci-dessous.

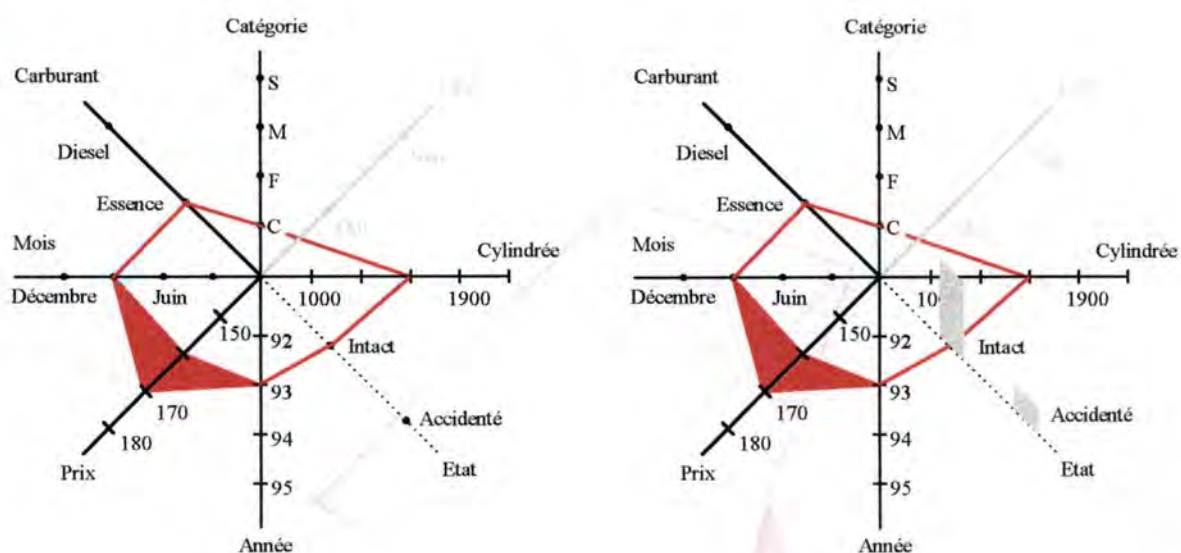


Figure 3.3. : En cliquant sur l'axe Etat, un graphique 2D des pondérations est affiché (passage de la partie gauche à droite de la figure).

Cette méthode de représentation a deux avantages majeurs :

1. elle fournit une image synthétique d'un objet symbolique complexe,
2. elle permet de comparer visuellement plusieurs objets entre eux.

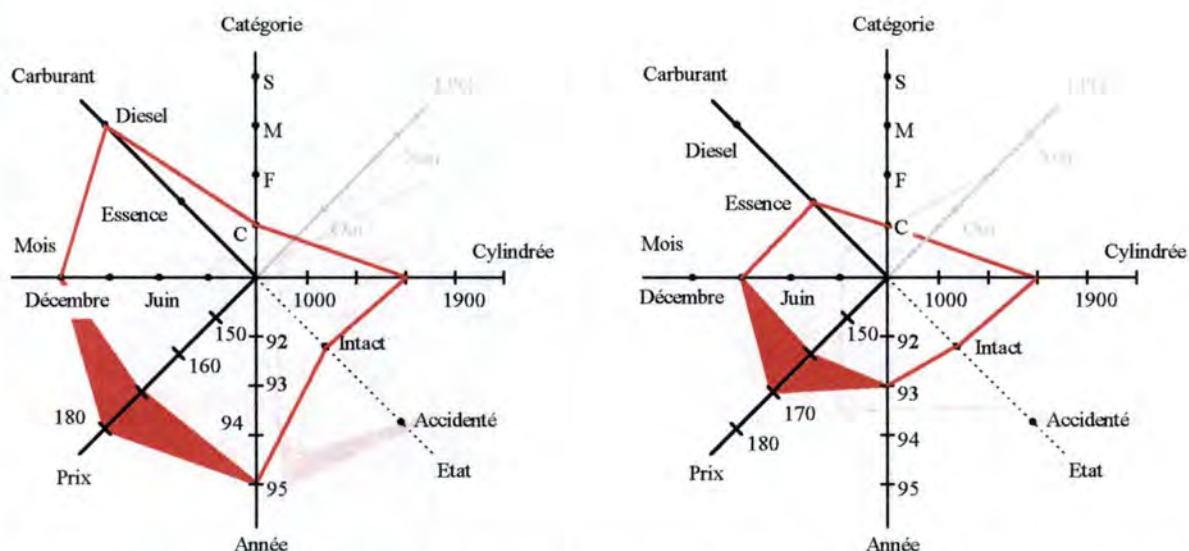


Figure 3.4. : comparaison de deux voitures représentées selon la méthode de l'étoile zoom.

Pour permettre une comparaison rigoureuse, il est nécessaire :

- de choisir un ensemble d'axes communs de référence, même s'ils ne sont pas « utiles » (compte tenu des valeurs manquantes) à tous les objets;
- d'utiliser les mêmes échelles pour un même axe, d'un objet à l'autre;
- de conserver le même ordre de présentation des axes.

### 2.3. Plus avant dans la représentation

Nous avons déjà évoqué la possibilité d'obtenir sur un axe supplémentaire la distribution des valeurs pondérées. Ces axes sont représentés en pointillés et, en cliquant dessus, la distribution est représentée.

Les dépendances entre variables et les taxonomies au sein de catégories ne sont pas automatiquement insérées sur le graphique.

De petits axes attachés aux catégories sont dessinés pour indiquer qu'une variable dépend de la valeur de cette catégorie.

La présence d'une taxonomie est indiquée par un icône représentant une hiérarchie, icône situé à proximité du nom de l'axe. La taxonomie correspond à une organisation des catégories d'une variable en une hiérarchie. Exemple : une variable TYPE DE MOTEUR possède les catégories suivantes : ELECTRICITE, THERMIQUE, SUPER, SANS PLOMB, DIESEL, GAZ. Les catégories peuvent être organisées de la manière suivante :

TYPE DE MOTEUR

↳ ELECTRIQUE

↳ THERMIQUE

↳ CARBURANT

↳ DIESEL

↳ SUPER

↳ SANS PLOMB

↳ GAZ



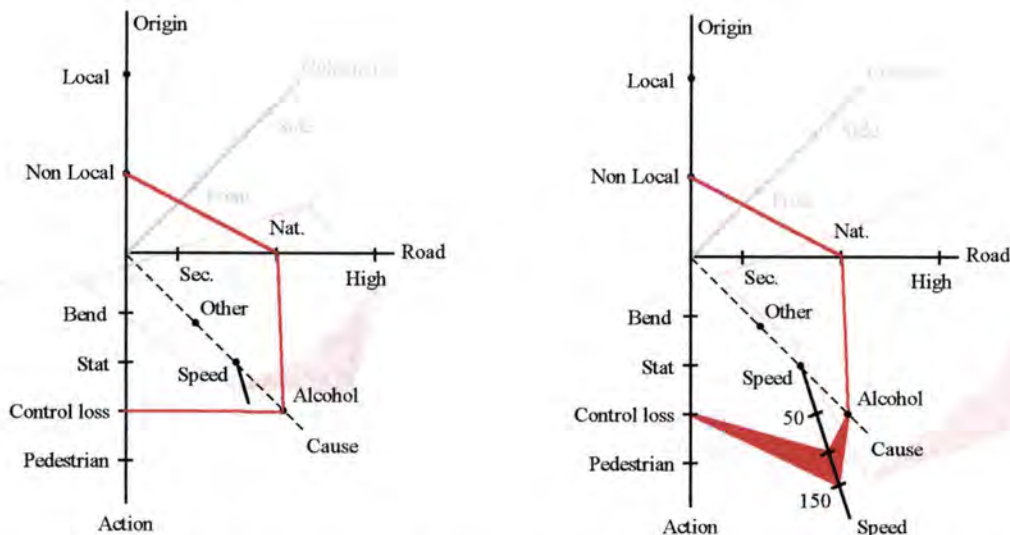


Figure 3.5. : Représentation de la dépendance entre catégorie d'une même variable, exemple tiré de [NOIRHOMME-FRAITURE 97b] et modifié.

Si l'utilisateur clique sur le petit axe (Speed de la variable Cause), la variable dépendante est ajoutée au graphique (de la partie gauche à droite sur la figure 3.5.).

Les auteurs permettent de représenter les représentations 2D et 3D, cela se justifie. Selon eux, les avantages des deux méthodes sont complémentaires :

- la représentation 2D est plus classique, habituelle, certaines personnes ont du mal à percevoir la 3D sur un écran,
- la difficulté d'identification des différences entre 2 objets lorsqu'ils sont représentés tous les deux en 3D,
- la 3D permet de représenter plus d'informations.

Le prototype tel qu'il est développé permet une animation du graphique en déplaçant horizontalement et/ou verticalement son point de vue (3D). Il est possible de représenter une étoile zoom 2D en 3D en y ajoutant verticalement les histogrammes correspondants aux variables probabilistes.

Lorsqu'une variable qualitative peut prendre plusieurs valeurs simultanément, la représentation doit pouvoir en tenir compte. Les différentes valeurs sont reliées aux axes voisins. A la figure 3.6. ci-dessous, l'objet représenté correspond à une voiture qui serait à la fois citadine et sportive. Par rapport aux intervalles, il ne s'agit pas de colorer la surface comprises dans les limites de l'étoile.

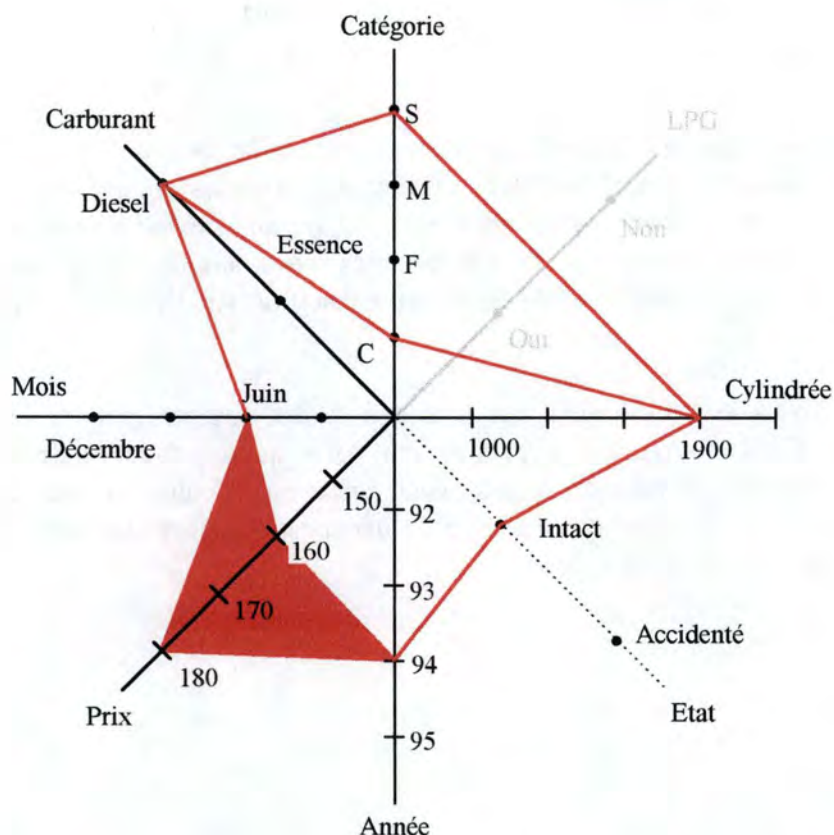


Figure 3.6. l'objet représenté correspond à une voiture qui serait à la fois citadine et sportive.

## 2.4. L'implémentation

Actuellement, les fonctions que nous avons présentées sont implémentées dans un environnement Windows. A court terme, le prototype inclura des fonctionnalités d'édition afin de permettre la modification directe des données.

Les données sont placées dans des fichiers « txt » classiques, ces données peuvent correspondre soit à la description d'un objet symbolique, soit à la description d'une instance d'un objet symbolique. Dans le premier fichier, les différentes caractéristiques d'un objet symbolique sont spécifiées, et dans le second, on indique les valeurs qui correspondent à une instance donnée de cet objet. Le logiciel permet, via une fonction, d'afficher ces données.



### 3. Conclusions

La représentation en 'Etoile Zoom' développée par M. NOIRHOMME-FRAITURE et M. ROUARD donne une image synthétique d'objets multivariés complexes. La technique est l'objet de tests par des utilisateurs appartenant à des organismes de statistique officiels de la CEE. A terme, d'autres possibilités et fonctionnalités seront ajoutées : l'apparition des données chiffrées dans un mini tableau modifiable, la représentation sur l'étoile de l'interaction entre deux variables, deux axes de l'étoile.

Nous l'avons vu, ces graphiques sont des étoiles compte tenu de la représentation d'axes radiaires. Elles sont qualifiées par « zoom » parce que, en fonction de ses besoins et de ses desiderata, l'utilisateur peut en sélectionnant un axe particulier obtenir des informations plus précises sur la variable en question : dépendance entre variables, représentation supplémentaire des histogrammes de poids.

### 1. Introduction

L'objectif d'une interface efficace est de permettre une visualisation maximale et rapide, clarté et convivialité s'imposent. C'est l'objet de la conception d'interfaces : édicter des règles ergonomiques pour atteindre l'objectif, règles à respecter dans la phase de développement de l'application.

A partir de 938 références bibliographiques, J. VANDERDONCKT dans [VANDERDONCKT 94] synthétise ces règles. Dans le cadre du projet FIRST, il a tenté de répondre à la question : « *Est-il possible d'avoir à sa disposition un volume unique et compréhensible rassemblant toute cette masse d'informations, de règles, de standards, de recommandations, ...?* »

Au-delà de cette compilation, l'auteur a développé un modèle E/A dont le but est de classer, rassembler et reformuler les règles ergonomiques. Nous renvoyons le lecteur intéressé aux publications de l'auteur : [VANDERDONCKT 94], [VANDERDONCKT 96a], [VANDERDONCKT 96b]. Dans le cadre du mémoire, nous ne pouvons que recommander de respecter les règles ergonomiques identifiées. Nous ne les reprenons pas telles quelles. Dans le cadre des trois chapitres qui suivent, nous tâchons d'émettre une série de recommandations pour l'utilisation du son, de la couleur et de la 3D au sein d'applications de représentation d'étoiles zoom.

J. VANDERDONCKT dans [VANDERDONCKT 94] justifie l'emploi de dessins. Utiliser le dessin, l'illustration pour :

- **motiver**, attirer l'attention de l'utilisateur (un dessin est plus agréable à visualiser);
- **expliquer** et décrire;
- **souligner** les points importants;
- **montrer** des interrelations complexes.

Ceci rejoint pleinement les objectifs qui ont motivé l'élaboration d'une méthode de représentation originale pour les objets symboliques.

Dans le cadre de ce chapitre, nous mettons en évidence les grands principes de la conception et de l'interprétation des interfaces. Notre objectif n'est pas l'exhaustivité, mais de jeter les bases sur lesquelles reposent les trois chapitres suivants.

### 2. Concevoir une interface

Compte tenu du caractère récent de l'utilisation de documents hypermédias digitaux et compte tenu de la difficulté conceptuelle d'intégrer plusieurs formes de médias au sein de présentations cohérentes, il n'y a pas de standards largement reconnus pour l'organisation visuelle de documents électroniques [LYNCH 94a, LYNCH 94b et VANDERDONCKT 96].



## **2.1. Les règles de conception**

A l'image de ce que réalise J. VANDERDONKT dans le cadre de son cours d'interface homme/machine, MULLET dans [MULLET 96] a tenté de synthétiser un ensemble exhaustif de principes de base de conception. Nous les résumons en les regroupant en six groupes.

### **2.1.1. L'élégance et la simplicité**

Si l'élégance apparaît immédiatement dans le succès d'une architecture graphique, sa simplicité est parfois exagérée. En fait, la simplicité d'une solution élégante est souvent son aspect le plus déroutant. Une solution simple révèle une compréhension profonde du problème et assure que les aspects essentiels soient saisis rapidement. Ce double aspect paradoxal (élégance et simplicité) au niveau de la quantité d'informations à traiter, doit trouver un équilibre :

- simplifier, affiner les éléments lors de la conception,
- combiner les éléments pour maximiser leur pouvoir d'expression.

Les conceptions les plus intéressantes sont souvent le fruit d'un processus continu de simplifications et d'affinages.

### **2.1.2. L'échelle, le contraste et les proportions**

La combinaison entre les trois domaines souligne l'harmonie d'une représentation. Lorsqu'un seul élément est trop large, trop étroit, trop petit, trop discret, trop sombre, c'est l'ensemble de la présentation qui en souffre. Pour maîtriser les relations entre échelle et contraste efficacement, il faut :

- intégrer figure, fond et champ de vision,
- mettre en évidence, souligner les limites d'éléments,
- établir une disposition en couches (superposées).

C'est ce domaine de la conception visuelle qui apparaît comme le plus difficile et celui qui requiert le plus de pratique. Les règles ne suffisent pas, l'expérience est nécessaire.

### **2.1.3. L'organisation de la structure visuelle**

C'est le premier aspect qui est perçu et qui va guider l'utilisateur au cours de l'interaction, il faut :

- assurer un équilibre en utilisant la symétrie,
- utiliser les alignements pour établir des relations visuelles,
- modeler l'affichage avec des espaces blancs (des espaces vides d'informations).

Sans l'intégrité offerte par une organisation visuelle cohérente et une structure logique, une architecture d'affichage peut rapidement devenir impossible à interpréter, à comprendre.



#### **2.1.4. Les modules en programmation**

Le concepteur doit établir un programme avec suffisamment de flexibilité pour permettre d'accommoder les demandes de chaque utilisateur.

Cette flexibilité doit être respectée tant que les éléments qui lui sont nécessaires, et donc présents en même temps, préservent la cohérence de la structure, pour cela il faut :

- renforcer la structure avec des répétitions,
- établir des unités modulaires,
- développer des schémas de construction de programmes basés sur des grilles existantes.

#### **2.1.5. Les images (illustrations) et représentations**

Toutes les considérations des points (1) à (4) s'adressent aux illustrations imagées (simplicité, ...). L'intérêt de l'image est notoire dans les situations pour lesquelles l'information qu'elle contient peut être comprise immédiatement et sans effort. Il faut :

- établir et maintenir une consistance visuelle,
- utiliser le bon format d'affichage,
- raffiner une image vers une abstraction progressive (résumer de plus en plus).

#### **2.1.6. Le style**

Le style aide l'utilisateur à imaginer un design en fournissant un contexte expérimenté dans lequel le public potentiel peut interpréter les signes de l'interface, d'où nécessité de :

- maîtriser le style,
- travailler avec plusieurs styles,
- étendre et développer le style.

MULLET 96 conclut en soutenant qu'une approche rationnelle de la conception d'interfaces visuelles est non seulement possible mais essentielle pour des applications graphiques. Que des règles soient formulées n'implique pas qu'elles doivent toutes être respectées en même temps. Mais une entorse ne doit s'appliquer à plus d'une règle à la fois, sauf peut-être pour des praticiens expérimentés.

Il est utile de se souvenir de ces principes lors de la conception de représentations d'étoiles zoom. La structure d'une étoile zoom se fige lorsque l'on détermine le nombre d'axes. De la simplicité, nous retiendrons qu'il faut limiter les surcharges textuelles au minimum (nom des axes, des graduations, ...). Les combinaisons des éléments de graphiques (axes, graduations, ...) sont relativement bien déterminées et conventionnelles (orientation des axes, règles de graduation, ...).

L'harmonie de la représentation est un facteur très important de la représentation des graphiques. C'est d'autant plus frappant si l'on compare plusieurs étoiles : les règles d'échelle, de contraste et de proportions doivent être les mêmes quelle que soit la représentation.



A priori, les représentations s'adressent à des experts, au-delà du respect des conventions de représentation (connues d'eux et donc familières), il est impératif de veiller à l'uniformisation des règles de représentation. Afin d'assurer un équilibre visuel, les axes doivent être répartis équitablement sur 360°. Dans un même domaine d'étude, d'une représentation à l'autre, la même structuration est conservée (même disposition des axes, des variables sur les axes, ...).

Les signes, abréviations, dénominations, le formalisme utilisé doivent être compris de tous. Le niveau de complexité de ces signes peut être élevé puisque c'est à un public expert que s'adressent les représentations, mais il ne faut pas que l'effort mental de compréhension, du formalisme, soit tel qu'il décourage l'effort d'interprétation de la représentation.

### 3. Manipuler les objets d'une interface

Pour manipuler les étoiles zoom représentant des objets symboliques et leurs composantes, différents moyens d'interaction<sup>1</sup> sont utilisables. Nous reprenons deux tableaux de synthèse réalisés par J. VANDERDONCKT et extraits de [VANDERDONCKT 94].

Légende du premier tableau :

C = clavier; E = écran tactile;  
F = flèche de déplacement; M = manette; S = souris.

Type de tâche	Utilisateurs novices	Utilisateurs experts
Sélection d'une petite zone	E et C plus rapides que F, S, M	C plus rapide que F
	C plus facile que S, M, F	Pas de différence
Sélection d'une grande zone	E plus rapide que C, F, S et M	E et C plus rapides que F
	E plus facile que F	E et C plus faciles que F
Manipulation d'objets	E plus rapide que M	E plus rapide que F
	Pas de différence	E plus facile que F

Tableau 4.1. : détermination des moyens d'interaction en fonction du type de tâche à accomplir, tableau tiré de [VANDERDONCKT 94] et modifié.

<sup>1</sup> Définition : chaque moyen d'interaction constitue un dispositif physique à l'aide duquel l'utilisateur acquiert et/ou restitue des informations relatives à sa tâche interactive.



Moyen d'interaction	Actions possibles	Inconvénients	Recommandé pour	Proscrit pour	Commentaires
Ecran tactile	Sélection	Activation accidentelle; fatigue des bras	Usages peu fréquents; pointage grossier	Usage continu; pointage précis.	Montage de l'écran pour fournir un repos des bras
Souris	Pointage; sélection; dessin; faire glisser; déplacement du curseur	Requiert de l'espace sur le bureau; possède un fil de raccord	Tâches requérant peu d'utilisation du clavier	Tâches requérant un changement fréquent entre le clavier et la souris	Peut intégrer des boutons fonctionnels avec le curseur
Styler optique	Déplacement du curseur, sélection, dessin	Souffre de la parallaxe, fatigue le bras	Utilisations peu fréquentes, tâches avec peu d'usage du clavier	Passages fréquents du clavier au styler optique et réciproquement; utilisation prolongée	L'écran doit être monté de façon à fatiguer le moins possible le bras en déplacement vertical
Manette	Suivi; sélection, déplacement du curseur	La souris peut être plus rapide pour la sélection de texte	Tâches comportant un positionnement intensif du curseur	Passages fréquents de la manette au clavier et réciproquement	Doit être disponible aussi bien pour les gauchers que les droitiers
Suiveur	Suivi, sélection, déplacement du curseur	La souris peut être plus rapide pour la sélection de texte	Intégration de graphiques avec saisies au clavier	Passages fréquents du suiveur au clavier et réciproquement	Doit être disponible aussi bien pour les gauchers que pour les droitiers
Clavier alphanumérique	Sélection, saisie textuelle, saisie numérique	Sous - optimale pour l'interaction iconique et la manipulation directe	Saisie de données à usage général	La sélection est plus lente par la frappe de la touche que par le pointage	Utilise une disposition standard identique à celles des machines à écrire.

Tableau 4.2. : avantages et inconvénients de quelques moyens d'interaction, tableau tiré de [VANDERDONKT 94] et modifié.

L'emploi de moyens d'interaction moins classiques (écran tactile ou la manette), peut-il faciliter l'analyse et la manipulation d'objets symboliques ?

C'est sans doute la phase de l'apprentissage d'utilisation qui pose problème. Puisqu'il faut apprendre à utiliser un nouveau logiciel, pourquoi ne pas utiliser de nouveaux moyens d'interaction en même temps ?

Pour manipuler des objets, l'écran tactile faciliterait la sélection des axes sur le graphique. Pour sélectionner et demander la visualisation de l'axe d'une variable probabiliste. Néanmoins, nous pensons que la souris suffit amplement.

La manette trouve éventuellement son intérêt pour la manipulation des objets 3D et pour permettre les rotations et déplacements de l'objet. Mais la souris peut, ici également, très bien remplacer ce moyen d'interaction. Le logiciel tel qu'il est prévu dans le futur s'adresse à



des utilisateurs experts, l'insertion de nouveaux moyens d'interaction (autres que la souris, le clavier et éventuellement le pavé numérique, outils familiers), obligerait des efforts d'apprentissage coûteux en temps. Nous ne pensons pas que cet investissement soit « rentable ».

En conclusion, l'utilisation de nouveaux moyens d'interaction peut, certes faciliter la manipulation et de là l'analyse visuelle des objets symboliques. Mais ces moyens doivent être des « plus » par rapport à ce qui existe habituellement (clavier et souris). Différents auteurs soulignent dans leurs publications l'emploi d'une souris spéciale permettant des déplacements dans les trois dimensions. En accord avec leurs constatations que nous synthétisons comme suit : le temps d'adaptation est souvent très important, l'effort d'apprentissage conséquent et les résultats ne sont pas à la hauteur des attentes, compte tenu de l'investissement en temps d'apprentissage et comparé à la souris habituelle. Néanmoins, ces moyens représentent des innovations et renferment des potentialités, nous en reparlerons.

## **4. Interpréter les éléments d'une interface**

La maîtrise du fonctionnement du système visuel est importante pour la conception d'interfaces visuelles. Par la représentation graphique, les données sont codées et le système visuel est chargé de décoder pour interpréter. C'est pourquoi la connaissance des caractéristiques de la perception humaine est fondamentale.

C'est une tâche très délicate que de vouloir assimiler entièrement les principes de la perception, parce qu'elle implique une compréhension en profondeur de la conscience humaine. Il n'existe pas de théories expliquant complètement tous les phénomènes de perception.

La théorie de MARR nous éclaire sur un point essentiel dans notre étude : la perception des formes [UNRUH 96]. La théorie neurophysiologique nous indique, elle, comment nous détectons les caractéristiques d'une image [ROSE 96]. Pour plus d'informations concernant les développements théoriques, nous renvoyons le lecteur intéressé aux références bibliographiques.

### **4.1. La perception des formes**

L'hypothèse de MARR définit qu'une représentation stable et fiable de la forme d'une image ne peut être dérivée en une seule étape. Un processus de perception de la forme en plusieurs étapes est proposé comme suit :

- 1 - *L'image rétinienne* : la distribution spatiale de l'intensité lumineuse sur la rétine est le point de départ du processus. Elle constitue le déclencheur de la procédure.
- 2 - *Le sketch primaire* : à cette étape, la série de données sur l'intensité est prise en considération et certains types d'informations sont extraits de celles-ci. Ce sont les premières



entités géométriques qui sont identifiées, la détection de surfaces devient possible. A ce stade l'approche computationnelle<sup>2</sup> intervient, dès la construction du sketch primaire.

- 3 - *Le sketch 2  $\frac{1}{2}$  D* : les détails sur la profondeur et l'orientation des surfaces visibles sont extraits, l'image commence à apparaître.

- 4 - *Le modèle de représentation 3D* : les surfaces, leurs positions, leurs orientations dans l'espace 3D deviennent explicites et sont reliées pour représenter des objets entiers. La représentation totale est maintenant décrite en termes d'objets. Le modèle du monde extérieur est achevé.

## **4.2. La perception des images**

Les figures, les images sont plus aisées à mémoriser que le texte et donc, plus fréquemment utilisées lorsque les informations doivent être retenues. Avant de représenter une image, conceptuellement il est nécessaire d'identifier ce que l'on veut représenter. La figure est-elle décorative, transformative, représente-t-elle des concepts ? Doit-on facilement identifier, interpréter ses composants ? Telles sont les questions auxquelles le concepteur doit répondre avant de « s'attaquer » à l'imagination, à la conception de la représentation.

L'utilisateur a souvent besoin d'instructions pour savoir comment regarder, interpréter une image. Sans lignes directrices de la part du concepteur à destination de l'utilisateur, la représentation pourrait ne pas remplir son rôle. Les instructions fournies peuvent apparaître clairement sous forme de liste de recommandations ou alors, la structure, le contraste sont tels que l'utilisateur est guidé inconsciemment. Si l'utilisateur expert s'attaque à une tâche familière, il a son propre schéma d'interprétation, il se l'est forgé au cours de précédentes analyses.

## **4.3. Le contraste<sup>3</sup>**

SHERRY souligne dans [SHERRY 95], l'attention est attirée par les zones de la représentation qui contrastent le plus de ses voisines. Les changements abrupts, nets sont les plus remarquables.

On gagne l'attention de l'utilisateur par un bit coloré au sein d'un écran noir et blanc, par un style d'écriture gras d'une police de caractères, par l'animation d'un objet, par des contrastes marqués dans la séquence des notes, du volume d'un extrait musical.

---

<sup>2</sup> Comme le laisse sous-entendre les termes, l'approche computationnelle traite la perception visuelle comme une procédure (procédure au sens informatique du terme), c'est-à-dire une suite finie d'opérations « programmables ». L'input de cette procédure correspond à un ensemble de deux images formées sur la rétine. L'output correspondra à une description symbolique du monde extérieur. Symbolique n'est pas employé dans le sens d'objets symboliques du premier chapitre mais veut dire que les objets du monde qui nous entoure sont représentés par des symboles qui se structurent entre eux pour former la représentation que nous nous en faisons.

<sup>3</sup> Le contraste : opposition de deux choses qui sont mises en valeur par leur juxtaposition (Larousse).



#### **4.4. La perception des graphiques, diagrammes**

Nous nous intéressons à la représentation en étoile d'objets symboliques; attardons-nous sur les recommandations que souligne SHERRY dans [SHERRY 95]. Les diagrammes et les graphiques illustrent des concepts abstraits. Les interrelations entre les éléments sont fortement conventionnalisées.

L'efficacité et la précision d'un graphique, d'un diagramme, dépendent de la localisation des composants et de leurs relations. Si plusieurs diagrammes sont utilisés en même temps, ou en séquence, c'est la même logique de représentation qui régit l'organisation de chaque illustration.

Augmenter la dissimilarité entre les éléments ne favorise pas la mise en évidence des relations. Tout dépend de ce que l'on veut montrer : la structure, les relations ou un élément face aux autres.

#### **4.5. La perception des sons**

Les sons sont organisés temporellement alors que les textes, les illustrations le sont de manière spatiale. Bien que ce soit le concepteur qui définisse les paramètres d'une séquence sonore, l'utilisateur doit être en mesure d'en contrôler les principaux (par exemple : si l'utilisateur écoute une série de directives qui doivent être précisément respectées). Le son représente le plus souvent l'inverse, le son est plus éphémère. Le texte, la représentation graphique sont plus efficaces pour transmettre de l'information complexe (on imagine mal de présenter un objet symbolique représenté par une étoile zoom sans un support graphique). Il faut bien évidemment que le son ait une parfaite concordance entre ce que contient le message verbal ou non et le contenu de la représentation graphique. C'est partiellement le cas lorsque des explications supplémentaires sont nécessaires à la compréhension d'une illustration.

SHERRY dans [SHERRY 95] montre qu'utiliser les séquences parlées est une des façons les plus expressives pour transmettre des instructions et des informations portant sur une représentation. C'est parce que le discours peut faire appel aux émotions, aux sentiments, aux sensibilités, que l'auteur soutient sa thèse. Nous verrons dans le chapitre suivant portant sur le son que les séquences sonores musicales peuvent, elles aussi, être des vecteurs informatifs performants.

### **5. Conclusions**

Dans l'attente de règles strictes de structuration des interfaces, certains auteurs ont identifiés des grands principes de conception d'interfaces. Que l'on aille du général (principe de MULLET) au particulier (règles ergonomiques de J. VANDERDONCKT), des lignes directrices, des guides existent pour aider le concepteur dans sa tâche de création, de structuration. Ces règles font souvent appel au bon sens et à l'expérience.

La représentation en étoile zoom, telle quelle est présentée, respecte assez bien les principes que MULLET décrit. Les moyens classiques d'interaction (souris et clavier) semblent suffisants pour la manipulation des étoiles zoom. S'adressant à des utilisateurs experts, la manipulation de différents éléments de la représentation peut être assurée par les moyens classiques que sont la souris et le clavier alphanumérique. Les périodes d'apprentissage et les efforts de familiarisation à de nouveaux outils sont importants pour ce type d'utilisateurs habitués aux moyens classiques.



### 1. Introduction

L'emploi du son, vocal ou non, au sein d'applications très interactives tels que les jeux vidéo, est devenu de plus en plus courant et populaire en raison de son énorme potentiel. Il peut être aussi utilisé pour présenter de l'information qui ne pourrait être l'objet de représentation via un mode de visualisation classique ou d'informations difficiles à discerner, telles que les données numériques multidimensionnelles. Le son est habituellement un complément aux *outputs* visuels parce qu'il augmente la quantité d'informations communiquées aux utilisateurs et/ou réduit la quantité d'informations que l'utilisateur doit traiter en mode visuel.

BREWSTER souligne dans [BREWSTER 93] que le système auditif est très puissant et trop souvent sous-utilisé dans les applications interactives. De plus, souligne le même auteur, il est évident que, psychologiquement, partager l'information entre plusieurs moyens de perception permet d'augmenter fortement les performances de l'utilisateur. Avoir de l'information redondante fournit à l'utilisateur deux chances d'identifier les informations, s'il ne peut se souvenir de la symbolisation d'une icône, il peut se rappeler des sons que cette icône véhicule.

A la différence de l'acuité visuelle, le son peut être entendu à 360° sans réel besoin d'une concentration sur un émetteur, ce qui fournit une plus grande flexibilité.

Autre différence : le son permet de capter l'attention de l'utilisateur pendant qu'il effectue une autre tâche (c'est le principe de la minuterie en cuisine). Sans jamais arriver à la même performance, seul le clignotement visuel permet d'attirer autant l'attention.

La plupart des objets que nous utilisons dans le monde réel produisent du son lors de leur manipulation. Même très faibles, nous les percevons; ils contribuent à construire notre représentation du monde.

Comme le souligne BEAUDOUIN-LAFON dans [BEAUDOUIN-LAFON 96], ces sons ne sont pas nécessairement naturels, l'essentiel est qu'ils convoient des informations utiles. Comme les illusions visuelles, les illusions sonores peuvent être insérées pour accentuer les *feed-back* des applications interactives.

Dans le cadre de ce mémoire nous ne proposons pas des exemples de sons utilisables dans telle ou telle situation. Nous invitons le lecteur intéressé par des suggestions de séquences sonores à consulter l'article de LEIMANN *et al.* [LEIMANN 95]. Dans le cadre de recherches, ces auteurs ont montré un ensemble de sons évocateurs, ces derniers ne s'adressent pas particulièrement à la compréhension de représentations d'objets symboliques.

Chaque utilisateur peut, à sa guise, insérer des sons et même des séquences vocales enregistrées. Des utilitaires, comme ceux des paramètres du panneau de configuration de WINDOWS 95, permettent d'insérer des sons enregistrés pour une longue série d'événements pré-définis. Nous appelons ces sons, des sons de signalisation. Nous développerons peu le sujet, si ce n'est pour évoquer des expériences montrant leur efficacité. Nous développerons plus en profondeur les sons comme moyens de perception, sons qui aident la compréhension de phénomènes.



## 2. Le son : moyen de signalisation

### 2.1. Introduction

Nous avons divisé ce chapitre selon le type d'informations véhiculé par le son. Utilisé soit comme moyen de signalisation, d'alerte, soit comme moyen de compréhension de phénomènes, le son aide à interpréter une figure, une représentation. Dans cette première partie, nous nous intéressons au son comme moyen de signalisation. BREWSTER, dans [BREWSTER 95], souligne un problème fréquent à l'emploi des boutons : le curseur de la souris glisse malencontreusement dessus et l'utilisateur provoque involontairement une action (lancement d'un programme, fermeture de fichier, ...). L'utilisation simultanée d'un son lié au bouton permet de réduire considérablement ce genre d'erreur (25 %) et diminue le temps de réaction de l'utilisateur face à certains événements (45 %).

### 2.2. Les expériences

DI GIANO *et al.* en 1992 dans [DI GIANO 92] ont ajouté du son à un environnement de programmation. Un son est attaché à un bloc de code ou de routine, d'où certaines erreurs de sémantique, comme les bouclages infinis, sont plus facilement détectées.


Les auteurs partent de l'hypothèse et de la question suivante : (1) même pour un programmeur, un logiciel peut être une boîte noire mystérieuse; (2) existe-t-il des outils, tel qu'un stéthoscope médical, pour écouter le rythme cardiaque d'un logiciel ?

La mise au point d'un environnement : *LogoMedia*, permet d'associer des zones musicales à des événements d'un programme en cours de développement. Ces associations aident à la compréhension et à l'analyse du comportement d'exécution du programme.

Les auteurs soulignent l'intérêt du son : pour eux, le son fait partie du *software visualization*. Grâce au son, se forment des images mentales du comportement, de la structure de logiciels par l'écoute des représentations acoustiques et de leurs spécifications ou exécutions.

En regard à certaines commandes, le programmeur décide d'associer un son. Une petite icône dans le cadre de programme signale la présence d'un son.

Les sons peuvent être placés, soit manuellement, soit de façon répétitive et automatique.

Par exemple :  : piano → Sum :a :b], signifie que quelques notes de piano sont jouées si la somme des variables a et b est modifiée.

L'idée d'un son typique pré-programmé peut être intéressante pour la mise en évidence automatique de certains événements (en particulier pour la représentation d'objets symboliques).

JACKSON et FRANCONI en 1994 dans [JACKSON 94] ont utilisé le son dans un environnement de programmation pour situer l'utilisateur au sein de l'exécution parallèle de programmes.



HAYWARD, en 1994, a démontré l'intérêt de l'utilisation du son sur des données brutes. En attachant des sons à des signaux sismiques : les événements importants sont rapidement mis en évidence par l'écoute des changements dans les données. Les possibilités d'écoute de données sont potentiellement infinies. Un exemple bien connu est celui de l'écoute des pulsations cardiaques en bloc opératoire.

L'emploi d'icônes sonores présente de grands intérêts et est de plus en plus expérimenté. MYNATT dans son article [MYNATT 94] propose une méthodologie de sélection, d'utilisation et d'évaluation d'icônes sonores. Ces derniers convoient de l'information symbolique vers l'interface auditive, tout comme les icônes graphiques vers l'interface visuelle. Nous synthétisons sous forme de recommandations la méthodologie de MYNATT.

Quels sont les facteurs déterminant l'emploi d'icônes sonores ?

- **Leur identification** : l'utilisateur doit pouvoir reconnaître un son sans équivoque, l'évaluation de son identification doit être réalisée préalablement auprès d'utilisateurs potentiels. Son apprentissage doit être le plus aisé possible.
- **Leur concordance conceptuelle** : le son doit coïncider avec l'aspect de l'interface que le concepteur veut indiquer à l'utilisateur. Le son doit transmettre chez l'utilisateur une image évoquant le message.
- **Leurs paramètres physiques** : l'emploi du son est influencé par ses paramètres physiques tels que l'intensité, la qualité, l'intervalle de fréquences. Le son doit être audible et « agréable ».
- **Les préférences de l'utilisateur** : la façon dont l'utilisateur réagit face à l'émission du son est importante : il faut éviter l'irritation par un bruit intempestif. La perception des bruits qui entourent est différente d'une personne à l'autre. Chacun a sa sensibilité. Il est utile de fournir à l'utilisateur la possibilité d'adapter les sons qui lui sont proposés. C'est à quoi s'est attaché GAVR. Dans [GAVR 93], l'auteur propose une série d'algorithmes permettant de modifier l'aspect sonore des *feed-back* qu'une interface met à la disposition des utilisateurs. Nous pensons que c'est un aspect très important et difficile à réaliser : que chaque utilisateur puisse dans un même environnement interpréter des messages en fonction de sa propre sensibilité.

Dans le reste du chapitre, nous nous intéressons à l'information visualisable. Nous voulons, ici, nous concentrer sur l'intérêt du son portant sur les objets cachés [BREWSTER 94] et [BLATTER 92].

Les auteurs, dans [BREWSTER 94], décrivent une méthode d'analyse des interactions pour identifier les situations où de l'information cachée existe et où des séquences sonores non vocales peuvent être utilisées pour dominer les problèmes de reconnaissance associés. Pour tester leur méthode, les chercheurs l'ont appliquée sur des ascenseurs de fenêtre.

En fonction du temps de réponse et du taux d'erreur, les auteurs affirment que l'utilisation de séquences sonores permet de réduire de façon statistiquement significative le temps nécessaire à l'exécution d'une tâche. De plus, soulignent-ils, l'effort mental nécessaire est nettement moins important.



Caractérisons d'abord les raisons pour lesquelles l'information peut être cachée :

- **L'information n'est pas disponible** : elle peut ne pas être disponible à cause des limites *hardware* de l'ordinateur, faute de puissance CPU, d'une taille d'écran trop petite.
- **L'accès à l'information est difficile** : l'information peut être affichable mais la difficulté repose ici sur les possibilités de l'obtenir. Par exemple [BREWSTER 94], pour connaître la taille et la date de création d'un fichier sur un MACINTOSH, il est nécessaire d'ouvrir une boîte de dialogue.
- **L'aire limitée du champ visuel** : l'information peut être cachée parce qu'elle est en dehors du champ visuel. L'utilisateur ne regarde pas au bon endroit, au bon moment, pour voir ce qui est représenté. Un des avantages du son est qu'il peut être entendu à 360° autour de nous. Les utilisateurs peuvent l'entendre sans pour autant se distraire de la tâche qui les occupe.
- **L'écran** : pour que des informations sur l'état du système soient visibles, cela implique une occupation de l'espace de l'écran. Un avantage du son, est de ne pas consommer d'espace écran, mais de l'espace mémoire.
- **Modes** : un mode est une situation où l'interprétation de l'information dépend du système dans lequel on se trouve. Par exemple, (SELLEN dans [BREWSTER 94]), une suite de caractères apparaissant sur l'écran d'un éditeur de texte, n'a pas la même interprétation si l'on sait qu'elle identifie une commande d'un langage de programmation utilisée dans un programme.

Décrivons, à présent les types de données cachées :

- **L'événement** : de durée très limitée, l'événement marque un fait important pour le système. L'événement peut être généré par l'utilisateur (un clic souris), ou par le système (l'arrivée d'un mail). Les événements peuvent être input ou output.
- **L'état** : toute information concernant l'état du système, état perceptible par l'utilisateur. L'état fait référence à un paramètre qui a toujours une valeur, tel qu'un indicateur. C'est un composant statique et continu dans le temps. Ces informations peuvent être affichées visuellement (par un graphique) ou auditivement ou par les deux modes de représentation combinés.
- **Le mode** (origine d'un événement) : différents modes peuvent provoquer le même événement et conduire à des effets différents. Si le mode n'est pas affiché, l'utilisateur ne soupçonnera pas les conséquences.

En fonction du type de donnée cachée, les auteurs [BREWSTER 94] et [BLATTER 92] déterminent les *feed-back* adéquats :

- **Les *feed-back* dépendent-ils des actions de l'utilisateur (du système) ?** Les événements sont des actions dépendantes; l'action de la part d'un utilisateur doit apparaître instantanément et durant un court instant. L'utilisateur clique sur la souris, c'est une action statique, le clic



sonore ne dure qu'un court instant. En ce qui concerne les *feed-back* sur les états, ils sont indépendants de l'activité, ils continuent dans le temps qu'il y ait de l'activité ou non.

- **Combien de temps doit durer le *feed-back* ?** Les événements sont passagers, ils apparaissent momentanément. Les informations concernant l'état sont continues dans le temps. L'émission sonore doit respecter ce facteur de continuité temporelle.

- **Le *feed-back* doit-il changer lorsqu'il est utilisé ?** L'information sur un état peut être statique ou dynamique. Par contre, les événements sont statiques, ils apparaissent un court moment et signale qu'une chose particulière se déroule.

- **L'utilisateur est-il obligé d'entendre le *feed-back* ?** Les événements sont souvent demandés car ils signalent des faits particuliers souvent importants. L'utilisateur prélève dans l'ensemble disponible d'événements signalables ceux qui l'intéressent. Il faut donner à l'utilisateur la possibilité de segmenter les événements, différencier ceux qui sont urgents des autres.

BREWSTER dans [BREWSTER 94] a testé sa classification et sa méthode sur les ascenseurs de défilement de fenêtre. Des *earcons* signalant des événements et des informations cachées sont utilisés. Leur emploi a permis de réduire de façon significative le temps de réaction des sujets cobayes. La combinaison des deux modes de perception, combinaison naturelle, permet d'accroître les performances en terme de temps de réponse et d'effort mental à fournir.

BREWSTER, dans [BREWSTER 97], utilise conjointement le son et la représentation graphique pour améliorer l'interaction. Il décrit une expérience où le son est employé pour faciliter la manipulation de menus déroulants. Une difficulté survient lorsque l'utilisateur glisse malencontreusement d'un item alors que, avec le souris il tente de sélectionner l'item en question. Des séquences sonores sont attachées à différents événements (ascenseur en bout de course en bas, en haut, curseur à côté de l'ascenseur, ...). En fonction de la classification, BREWSTER détermine les *feed-back* les plus judicieux et les teste. Les résultats de l'expérience montrent une réduction significative de l'effort subjectif requis pour accomplir la tâche de sélection; ainsi qu'une réduction statistiquement significative du temps de réalisation de la tâche.

### **2.3. Conclusion**

La combinaison des informations graphiques et sonores sur une même interface est devenue naturelle. Chaque jour, les deux sens sont associés pour permettre la perception optimale d'informations complémentaires du monde. Les deux sens sont interdépendants.

Le système visuel nous donne des détails d'un foyer, là où le système auditif fournit des informations générales de ce qui nous entoure, nous alertant de choses que nous ne pouvons voir. Les deux sens combinés nous fournissent toute l'information (parfois plus) dont nous avons besoin pour comprendre l'environnement. Pourquoi ne pas utiliser ces avantages au sein d'interfaces interactives ?



### 3. Le son : moyen de communication d'informations

#### 3.1. Introduction

Au-delà de l'aide à la perception que sont les sons de signalisation que nous avons décrits au point précédent, ne peut-on pas aller plus loin ? Comment utiliser le son comme moyen de perception, de représentation et pas seulement comme moyen de signalisation ?

#### 3.2. Les expériences

##### 3.2.1. Son et 3D : l'expérience de MEREU

MEREU et KAZMAN dans [MEREU 96] ont décrit dans cet article publié en 1996, une étude concernant la compréhension de graphiques en trois dimensions à l'aide de l'information sonore. Nous résumons ici leurs recherches.

##### A. Introduction.

Les *feed-back* audios ont déjà prouvé leurs apports significatifs pour des interfaces homme-machine d'applications numériques.

Le son y est utilisé:

- lorsque la présentation de l'information ne peut être formulée autrement,
- lorsque l'attention de l'utilisateur doit être focalisée sur un sujet particulier, ou si l'utilisateur effectue une autre tâche simultanément. Cette seconde occupation peut très bien n'avoir aucun rapport avec la première.

Par exemple, une machine de distribution de boissons émet habituellement un bruit sous forme d'un rythme régulier. Lorsqu'une anomalie apparaît, le rythme change pour alerter l'utilisateur. L'idée du son dans les applications est d'autant plus efficace que la tâche à réaliser est complexe. Dans ce cas, l'aide est d'autant plus significative. Les recherches dans le domaine sont récentes mais nombreuses.

Un autre créneau intéressant des *feed-back* audios est de permettre l'accès aux applications implémentées aux défailants visuels.

MANSUR et BLATTNER en 1985 dans [MEREU 96] ont montré que le son permettait de percevoir la relation existant entre X et Y, deux variables représentées sur un graphique. Pour ce faire, l'axe des abscisses est balayé et la hauteur du son est ajustée en fonction de la valeur de Y.

EDWARDS en 1988 dans [MEREU 96] identifiait les fenêtres avec un son. L'action de déplacer le curseur sur une fenêtre, de cliquer dedans provoque l'émission d'un son identifiant.

Le son peut aider les personnes défailtantes visuellement pour des applications 2D mais aussi pour des applications 3D. Les auteurs s'interrogent également sur l'aide que peut présenter le son pour les personnes d'acuité visuelle normale.



## B. Preamble.

Avant d'étudier l'effet du son sur la perception de personnes déficientes visuelles, les auteurs ont étudié les effets sur les personnes voyantes.

Quatre axes de recherche sont identifiés:

- déterminer si le son peut être utilisé comme une information reflétant la profondeur pour la perception de l'espace 3D,
- comparer les différences entre une variété d'informations sonores,
- noter le comportement de l'utilisateur quand il ne dispose que d'informations sonores,
- étudier l'influence d'interfaces sonores sur les capacités d'apprentissage des utilisateurs.

Les chercheurs (MEREU et KAZMAN dans [MEREU 96]) mettent au point trois environnements sonores différents (tonal, musical et orchestral).

**L'environnement tonal** attache une simple onde sinusoïdale à chaque position. Comme le curseur bouge dans l'espace 3D, le son va changer dans trois de ses dimensions que sont par exemple : sa balance, sa hauteur, son volume. Pour ce premier environnement qualifié de tonal, la meilleure répartition pour les dimensions spatiales face au son déduite par les auteurs est:

X (gauche/droite) → balance (gauche/droite),  
Y (haut/bas) → hauteur (aigu/basse),  
Z (loin/près) → volume (calme/fort).

**L'environnement musical:** un morceau musical est joué, le morceau est modifié en fonction de la localisation du curseur. Le meilleur système identifié par les auteurs est:

X (gauche/droite) → balance (gauche/droite),  
Y (haut/bas) → hauteur (aigu/basse),  
Z (loin/près) → volume (calme/fort).

Dans un **environnement orchestral**, l'idée est de jouer un morceau de musique en arrangeant un groupe d'instruments dans une configuration standard et en l'utilisant pour localiser le son en 2D. Le but est d'apporter plus de précision quant à la localisation exacte dans ce plan. Le plan Z-X a été fractionné en 8 zones d'instruments différents comme le montre la figure 5.X :



Figure 5.1 : Arrangement de « l'orchestre », tiré de [MEREU 96].



Après de nombreux essais, les auteurs ont sélectionné les huit instruments les plus identifiables parmi les 128 disponibles dans le set *MIDI*.

En utilisant cet environnement orchestral, la musique jouée, correspondant à la position du curseur, est dominée par l'instrument de la section où le curseur se trouve. Cette amélioration (8 sections sur le plan de Z-X) est intégrée car beaucoup d'utilisateurs ont des difficultés à identifier les différences de hauteur du son.

### ***C. Résultats de l'expérience.***

Les résultats montrent que les trois environnements sonores (tonal, musical et orchestral) réduisent les erreurs de positionnement par rapport aux environnements n'utilisant pas le son : par contre, le temps nécessaire à la réalisation de la tâche est largement augmenté (cfr. le tableau 5.1.).

	Tonal	Musical	Orchestral
Erreurs de positionnement	78,9 %	61,1 %	32,4%
Temps de réalisation	123,4 %	178,4%	215,1 %

Tableau 5.1. : les erreurs de positionnement et le temps de réalisation des tâches dans les environnements sonores, tiré de [MEREU 96] et modifié.

Nous rappelons que cette expérience s'est déroulée avec uniquement des sujets « bien » voyants. Les résultats des expériences conduisent les chercheurs à déduire qu'il est possible d'utiliser des environnements sonores dans les applications utilisant exclusivement le son quand l'utilisateur (voyant) doit se focaliser sur une autre activité ou si l'utilisateur est mal-moins voyant. Le modèle environnemental qui semble le plus performant est le modèle tonal. La différence entre les caractéristiques des sujets, comme le sexe, le type de formation, ses capacités musicales ou l'expérience d'interfaces graphiques ne jouent pas de façon significative. Les auteurs concluent également, qu'avec un peu d'entraînement, le temps de réalisation est diminué : les utilisateurs s'habituent aux conditions.

### ***D. Les expériences avec des mal-moins voyants.***

Les trois mêmes environnements (tonal, musical et orchestral) sont testés dans les mêmes conditions de réalisation que pour les sujets voyants, hormis le fait qu'ici seul le son intervient.

### ***E. Les résultats.***

Les trois environnements sont comparés en utilisant une table ANOVA. Nous synthétisons les résultats à l'aide de graphiques et de commentaires explicatifs.



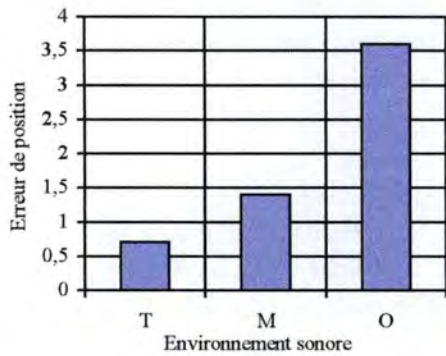


Figure 5.2. : Environnements sonores et erreurs de position, tiré de [MEREU 96].

La figure 5.3. montre qu'aucune différence significative n'apparaît dans le temps de réalisation de la tâche selon les trois modèles environnementaux.

Cette figure 5.2. montre le niveau d'erreur de positionnement suivant les trois modèles d'environnement. Le niveau d'erreur de l'environnement tonal est approximativement de moitié de celui de l'environnement musical et le cinquième de l'orchestral.

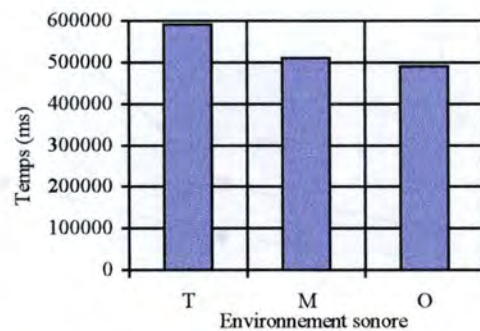


Figure 5.3. : Environnements et temps (ms) de réalisation des tâches, [MEREU 96].

Ce graphique 5.4. illustre l'utilisation de quatre paramètres subjectifs énoncés en fin d'expérience par les différents utilisateurs mal-moins voyants : la préférence, la performance, l'utilisabilité et la rentabilité de l'environnement.

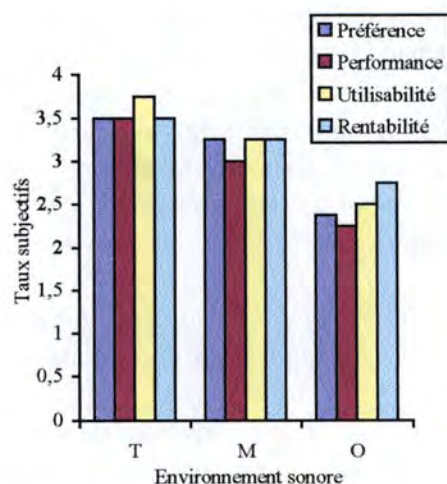


Figure 5.4. : appréciations subjectives dans les environnements sonores, [MEREU 96].



Les différences entre les trois environnements n'ont pas atteint un seuil qui témoignerait de différences statistiquement significatives entre eux. Ceci peut être une conséquence du faible échantillonnage. L'environnement tonal est le mieux classé par les individus déficients visuellement en fonction des 4 critères subjectifs.

*Comparons à présent les deux types de sujets.*

L'erreur de positionnement et le temps nécessaire à la réalisation d'une tâche entre sujets voyants et non voyants, montrent une interaction significative avec l'environnement sonore. Les figures 5.5. et 5.6. montrent respectivement ces interactions de l'erreur de position et du temps.

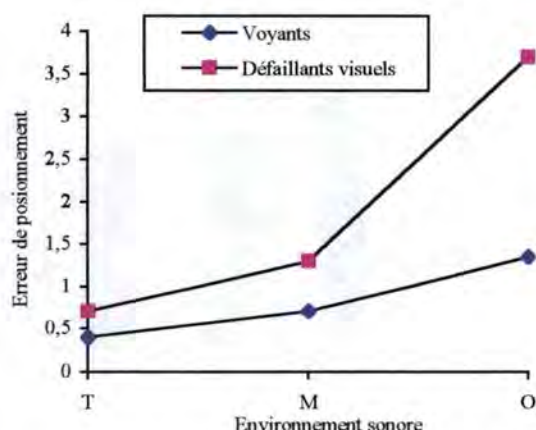


Figure 5.5. : Interaction entre la capacité visuelle face à l'erreur de positionnement, [MEREU 96].

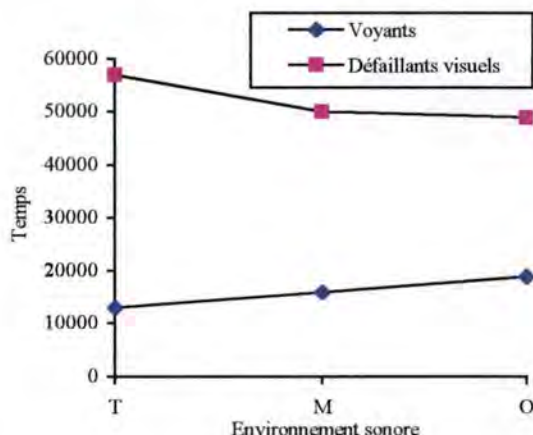


Figure 5.6. : Interaction entre la capacité visuelle face au temps (ms) de réalisation, [MEREU 96].

Il n'est pas surprenant de constater que les individus voyants ont une meilleure précision et un temps de réponse plus court puisqu'ils utilisent l'ouïe et la vue, alors que les personnes déficientes visuellement ne disposent que de l'information sonore. La différence de performance est significative. Il est intéressant de constater que les sujets non voyants prennent en moyenne 4,3 fois plus de temps en utilisant le modèle tonal. L'environnement tonal fournit, en gros, autant d'informations aux utilisateurs non voyants que n'apporte la combinaison d'informations visuelles et sonores pour les utilisateurs voyants. Ceux-ci positionnent très rapidement le curseur dans un plan X-Y, il faut par contre vraiment beaucoup de temps pour la localisation dans le troisième plan Z.

Intéressons nous maintenant à une expérience concernant les voyants avec uniquement des environnements sonores, l'acuité visuelle n'est plus mise à contribution.

Les deux graphiques ci-dessous (5.7. et 5.8.) sont de même allure que les précédents, mais ici les deux types de sujets (voyants et non voyants) ne travaillent que sur des applications sonores. Les individus déficients visuellement ont un nombre d'erreurs de positionnement significativement moindre que les sujets voyants et ce dans les environnements tonal et musical.

D'une manière générale, les sujets voyants prennent plus de temps et commettent plus d'erreurs que les sujets non voyants pour la réalisation d'une même tâche.



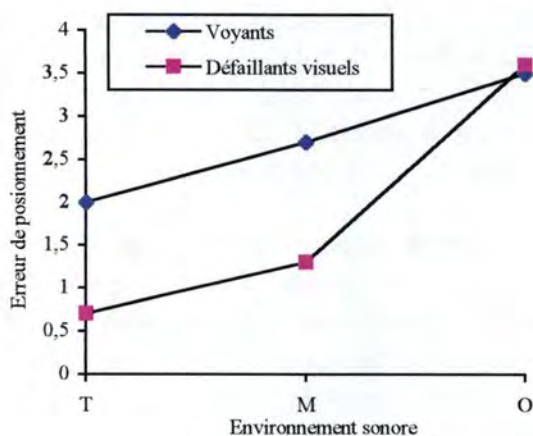


Figure 5.7. : Interaction entre la capacité visuelle face à l'erreur de positionnement, [MEREU 96].

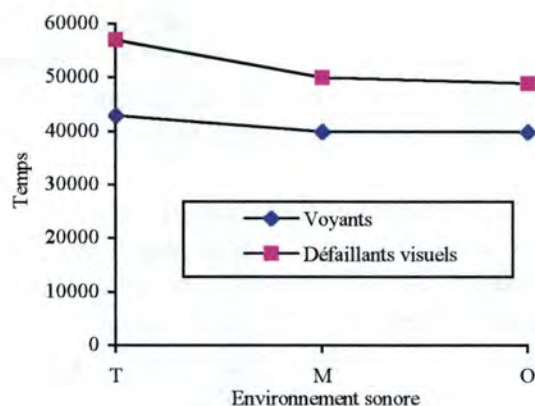


Figure 5.8. : Interaction entre la capacité visuelle face au temps (ms) de réalisation, [MEREU 96].

### *Quel est l'environnement faut-il préférer ?*

Des mesures subjectives réalisées (performance, préférence, convivialité et rentabilité) ont été comparées en utilisant le test ANOVA pour identifier d'éventuelles différences entre les goûts des voyants et ceux des mal-moins voyants.

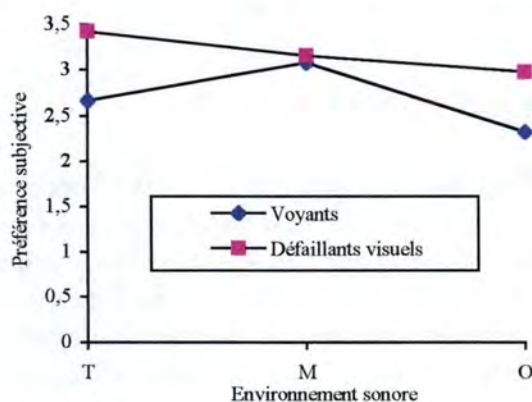


Figure 5.9. : Interaction entre la capacité visuelle face aux préférences d'utilisation, [MEREU 96].

Ce graphique 5.9. montre les différences concernant les environnements sonores étudiés faces aux préférences des individus testeurs des deux groupes. Les sujets défailants visuellement préfèrent l'environnement tonal alors que les voyants se tournent plus vers l'environnement musical.



## F. Conclusions.

Ces résultats ont montré que les sujets déficients visuellement peuvent utiliser les environnements sonores pour percevoir la profondeur et la position dans des applications 3D. Bien que pour eux le temps de localisation est significativement plus long que pour des personnes non atteintes, la précision de localisation est la même. Ceci est important pour l'intégration des non voyants dans le monde de l'utilisation de telles applications.

Subjectivement, les utilisateurs non voyants préfèrent l'environnement tonal alors que les voyants optent plutôt pour l'environnement musical. Cela dépend du rôle que le sujet octroie au son : pour le non voyant, le son représente sa seule source d'information, il n'est donc pas question d'alourdir la perception par des fioritures musicales. Par contre, pour les individus voyants, l'information sonore représente une seconde source et donc un environnement musical plaisant semble alors plus agréable.

L'environnement orchestral se classe mal dans les deux cas, deux explications contradictoires peuvent être avancées:

- le classement est d'autant plus valable que son le niveau de précision est accru,
- la conception moins familière d'une grille divisée en 8 zones augmente le degré de difficulté d'utilisation. L'effort à fournir en est d'autant plus important.

Remarquons que des *feed-back* audios permanents agacent les utilisateurs. Un environnement qui minimise ce problème est un environnement qui sera fréquemment utilisé.

En conclusion, le son peut aider les deux types de sujets (voyants et mal-moins voyants) dans leur utilisation des applications 3D. L'environnement sonore préféré dépend des capacités visuelles des utilisateurs: environnement tonal pour les personnes déficientes visuellement et environnement musical pour les autres.

### 3.2.2. L'utilisation d'*earcons* [BREWSTER 93]

Avec l'expérience de BREWSTER [BREWSTER 93], nous introduisons une notion sur laquelle beaucoup d'auteurs dissertent : les *earcons*. Les *earcons* sont, en résumé, des sons synthétiques qui peuvent être utilisés dans des combinaisons structurées pour créer des messages sonores qui représentant des parties d'une interface. BLATTNER dans [BREWSTER 93] définit les *earcons* comme des messages audios non verbaux utilisés en interface homme/machine, pour fournir à l'utilisateur, des informations concernant des objets, des opérations, des interactions.

Les *earcons* sont composés de séquences courtes et rythmées de notes. Les séquences ont une intensité variable, un timbre et un registre<sup>1</sup> qui peuvent également varier.

Dans leurs expériences, une évaluation des *earcons* est menée pour déterminer l'efficacité de communication du son. Une première expérience montre que les *earcons* sont plus rentables que l'émission de salves non structurées de sons et que le timbre musical est plus efficace que le son simple. Une seconde expérience est conduite pour améliorer quelques faiblesses mises en évidence dans le courant de la première expérimentation. Nous le verrons, les *earcons* sont des moyens efficaces de communication d'informations.

<sup>1</sup> Le registre en musique correspond à chacune des trois parties (grave, medium, aigu) qui composent l'échelle sonore ou la tessiture d'un son, d'une voix.



La première expérience est construite de façon à déterminer si les sons structurés tels que les *earcons* sont meilleurs que les sons non structurés pour véhiculer de l'information. De simples sons sont comparés à des timbres musicaux plus complexes. Rythme et notes sont aussi testés afin de différencier les *earcons* entre eux. Le rythme est une des méthodes les plus puissantes pour différencier les sources sonores. A la figure 5.10. nous reproduisons les rythmes et les notes utilisés au cours des quatre phases de l'expérience. Ces quatre phases sont séquentielles.



Figure 5.10. partitions utilisées dans le cadre de la première expérience, [BREWSTER 93].

Trois ensembles de sons sont construits pour la réalisation de l'expérience :

- 1) l'ensemble musical : des timbres musicaux différents produits par un synthétiseur (instruments différents);
- 2) l'ensemble simple : composé de timbre sonore simple (onde sinusoïdale);
- 3) l'ensemble de contrôle : juste composé de salves de sons informatiques classiques : des *beeps*. Contrairement aux deux premiers sets, il n'y a pas de rythmes qui interviennent ici.

#### Phases I :

Les sujets sont face à un écran (de type MACINTOSH) présentant différentes icônes de fichiers ou de programmes ainsi que des dossiers. Chaque famille d'objets (icônes d'un même programme, toutes les icônes de PAINTBRUSH, programme, documents, dossiers) partage le même timbre (le même instrument). Au sein d'une même famille les items sont différenciés par les notes jouées. Les sujets disposant de l'écran, font circuler la souris ce qui provoque l'émission d'une séquence en passant au-dessus d'une icône. Puis l'écran s'efface et les sujets doivent, par la seule écoute des séquences, déterminer l'information au sujet de l'application et du type d'item.

#### Phase II :

Le même type d'expérience est mené en relation avec des menus. Chaque menu a son propre timbre et chaque item de menu est différencié des autres par le rythme, les notes et l'intensité.

#### Phase III :

C'est une phase de re-test de la Phase I, mais cette fois sans période d'apprentissage. Les « questions » sont présentées dans un ordre différent. Le but est d'observer si les sujets peuvent mémoriser l'ensemble original d'*earcons* tout en ayant utilisé d'autres.

#### Phase IV :

C'est une combinaison des phases I et II, de nouveau sans période d'apprentissage.

La figure 5.11. ci-dessous présente les résultats.



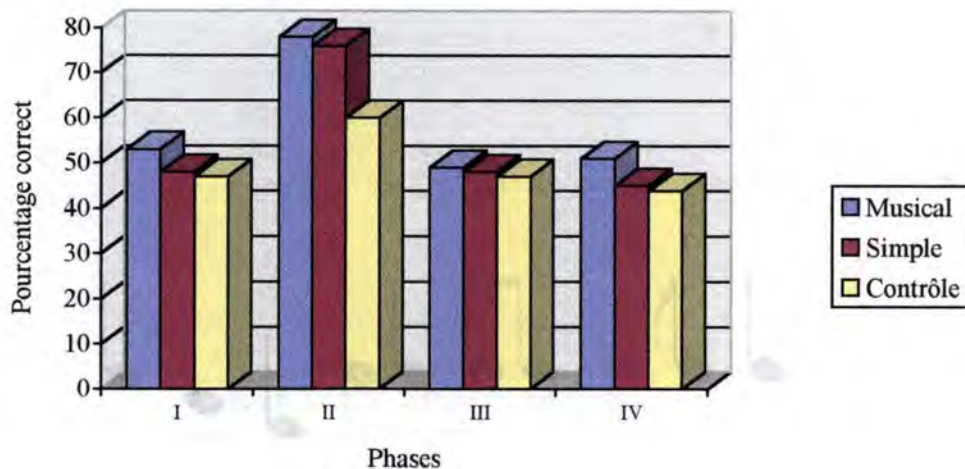


Figure 5.11. : Pourcentages de tous les scores répartis par phases, expérience 1, extrait de [BREWSTER 93] et modifié.

Même si les différences ne sont pas toujours statistiquement significatives, il apparaît que les *earcons* de type musical sont les meilleurs, quelle que soit la phase.

#### Phase I :

Cela indique que le timbre musical est plus facilement reconnu que les sons simples. L'utilisation du rythme ne permet pas d'améliorer nettement les performances; les différences ne sont pas statistiquement significatives.

#### Phase II :

Les différences sont plus significatives que pour la phase I. L'introduction de rythmes différents dans les séquences musicales augmente les performances. Le rythme s'il est correctement utilisé peut vraiment aider la reconnaissance. L'utilisation de note seule (de *beep*) est très difficile à interpréter.

#### Phase III :

Les scores ne sont pas significativement différents de ceux de la phase I. Cela signifie qu'après une courte période, les sujets sont capables de se souvenir des *earcons* entendus.

#### Phase IV :

Le fait de mélanger les menus, les items de menus, les icônes de fichiers, de programmes et de dossiers augmente la difficulté de reconnaissance. Les différences sont plus significatives entre les sets utilisés.

En général, sous certaines conditions, les *earcons* sont plus performants que les salves de notes non structurées ou que les bruits simples. Même si cela n'est pas toujours montré significativement, l'instrument joue un rôle important dans la reconnaissance. L'utilisation de rythmes différents mérite d'être étudiée de plus près.



Dans le cadre de la seconde expérience, les auteurs introduisent de nouveaux *earcons*, il y a, à présent, une plus grande différence entre eux. Comme le montre la figure 5.12. ci-dessous, et en comparaison de la figure 5.10. précédente, c'est en particulier sur le rythme que les auteurs jouent.



Figure 5.12. : partitions utilisées dans le cadre de la seconde expérience, [BREWSTER 93].

Le déroulement de l'expérience est similaire à celui de la première décrite ci-dessus. Seuls les sujets changent bien que leurs caractéristiques soient semblables pour permettre les comparaisons, et les *earcons* sont plus individualisables. L'analyse de la figure 5.13. permet de montrer les performances mesurées dans les deux expériences.

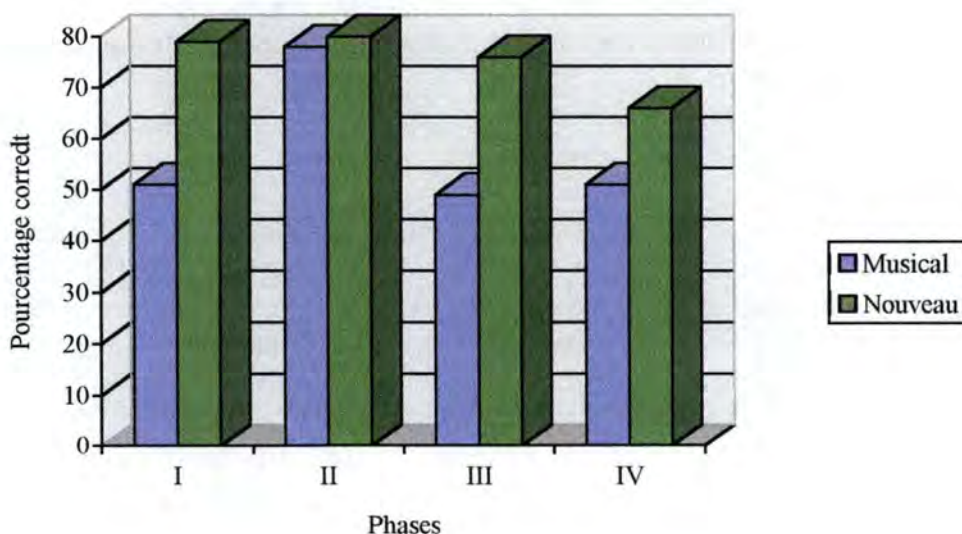


Figure 5.13. : Pourcentages de tous les scores répartis par phases, comparaison de l'expérience 1 et 2, extrait de [BREWSTER 93] et modifié.

En général, les performances sont améliorées. L'utilisation de nouvelles notes, de nouveaux rythmes, de registres plus différenciables permet de conclure que les séquences sont plus identifiables par les utilisateurs. La différence la plus notable se marque au niveau de la phase III, ce qui peut vouloir dire que la mémorisation est d'autant meilleure.

*Le fait que les utilisateurs soient musiciens ou non peut-il intervenir ?*

Dans le cadre de ces expériences, les auteurs ont sollicité des sujets musiciens et d'autres non musiciens. Quelle que soit l'expérience, les résultats de l'un ou l'autre groupe



n'étaient pas différents, les performances en reconnaissance étaient grosso modo identiques. Certains items de menus ou famille de programme sont parfois mieux identifiés par les individus qui ont l'oreille musicale. Forts de leurs expériences, les auteurs affirment que les performances de sujets ayant un bagage musical ou celles d'individus n'en ayant pas, sont identiques. Des différences plus importantes peuvent apparaître lorsque des changements plus subtils entre *earcons* sont utilisés. Dans ce cas, avoir une oreille musicale améliore ses performances.

### 3.3. Synthèse des résultats

Les *feed-back* audios ont prouvé leurs apports significatifs pour des interfaces homme-machine d'applications interactives. L'emploi du son est justifié lorsque la présentation de l'information ne peut être formulée autrement, ou lorsque l'attention de l'utilisateur doit être focalisée sur un sujet particulier, ou encore si celui-ci effectue une autre tâche en même temps.

Nous pouvons retenir de ces expériences que :

- les *earcons* facilitent la communication d'informations;
- l'utilisation de séquences de notes est préférable aux sons simples (simplistes);
- les capacités de mémorisation sont d'autant plus importantes que les séquences sont différentes entre elles;
- le fait d'être musicien n'augmente pas de façon significative les performances d'utilisation d'interfaces sonores.

L'expérience menée par MEREU est particulièrement enrichissante et étonnante : il est possible d'utiliser le son (voire de remplacer le mode de perception visuel par le mode auditif) pour permettre d'interpréter des représentations graphiques. Ces expériences montrent à quel point notre environnement sonore est fondamental et soulignent aussi la relative inexploitation actuelle de ces potentialités en informatique.

RIGAS et ALTY, dans leur expérience [RIGAS 97], examinent la possibilité de communiquer des informations de nature spatiale à des utilisateurs malvoyants. Le développement d'un outil (AUDIOGRAPH) permet aux utilisateurs déficients visuellement d'apprécier (d'interpréter) et de manipuler des objets contenant des zones graphiques. Une des phases de l'expérience correspond à une phase de création de représentations de formes géométriques simples et ce, par la seule information sonore. Il s'agit de dessiner des cercles des rectangles, .... Le son aide à la confection de formes géométriques; les résultats sont encourageants.

## 4. Conclusions

Des résultats des différentes expériences que nous relatons et à partir de l'étude de la littérature, nous tirons quelques lignes directrices dans l'utilisation de séquences musicales au sein d'applications fortement interactives.

Plus les différences entre les séquences présentes au sein d'une même application sont importantes, plus les performances de reconnaissance sont élevées.



Si les seuls changements relèvent de la subtilité ou sont de faibles amplitudes, ils ne seront remarqués que par les utilisateurs ayant développé des capacités musicales.

Nous synthétisons, pour chaque paramètre accessible à l'utilisateur, les principales recommandations d'emploi. A l'image de ce qu'ont réalisé MULLET ou J. VANDERDONCKT, il s'agit de conseils de niveau conceptuel. Du respect de ceux-ci dépendent la réussite et l'efficacité de l'interface. Le respect de ces règles n'est pas strict et absolu.

Concernant le **timbre** : BREWSTER suggère d'utiliser des synthétiseurs, la multiplicité des harmoniques y est possible. Cela aide à la perception et évite le masquage, la superposition. Des timbres différents ne doivent être mis en oeuvre que s'ils sont bien identifiables entre eux, par exemple : utiliser un instrument de la famille des cuivres et un piano plutôt que deux cuivres.

Les **notes** : il ne faut les utiliser que s'il y a de réelles différences entre elles. L'usage simultané de notes et de rythmes différents est très efficace. BREWSTER, de nouveau, propose d'utiliser au minimum 125 Hz à 150 Hz, c'est à dire une octave sous le Do et au maximum 5 kHz, ou encore 4 octaves au-dessus du Do.

Le **registre** : si l'on joue seulement sur ce facteur pour permettre la différenciation, alors il faut de grandes différences : au moins trois octaves.

Le **rythme** : le rythme est un très bon moyen de différenciation. Il n'est pas nécessaire d'avoir une oreille musicale pour bien identifier de faibles différences. Néanmoins, plus les différences sont importantes, plus les utilisateurs seront sensibles et performants. La combinaison de notes et de rythmes différents est le meilleur moyen d'identification. PATTERSON en 1982 dans [BREWSTER 93] souligne que pour des rythmes similaires, utilisant des notes différentes, les identifications de différentes séquences sonores demandent beaucoup d'énergie à l'utilisateur, ce qui va à l'encontre de l'objectif d'aide à la représentation. SUMIKAWA dans [BREWSTER 93] recommande l'emploi de séquences de 8 notes minimum. La longueur d'émission d'une note ne doit pas être inférieure à 0,125 sec.

L'**intensité** : PATTERSON dans [BREWSTER 93] suggère un intervalle de [10 dB, 20 dB]. L'intensité d'émission : trop faible, elle ne permettra pas l'identification, trop haute, elle gênera l'utilisateur. Ce paramètre doit être sous le contrôle de l'utilisateur. Il est clair que l'ensemble des séquences musicales ou sonores d'une même application doivent être émises dans la même zone d'intensité : l'utilisateur ne doit pas être obligé de jouer avec le volume ou être surpris en cours d'application. Si le son est trop fort, il peut ennuyer l'utilisateur et dominer les autres. S'il est trop faible, il peut être perdu, inutile, de même que l'information véhiculée.

Les **combinaisons** : tout en respectant les recommandations qui viennent d'être soulignées, c'est en combinant ces paramètres que les meilleures performances sont remarquées. Lorsque l'on utilise plusieurs séquences en même temps, il faut bien les séparer : un délai de 0,1 seconde est le minimum adéquat.



Nous l'avons vu au cours de l'expérience de BREWSTER, l'utilisation plus ou moins harmonieuse des notes musicales, plutôt que de salves de bruits, est plus performante. Réaliser des séquences de notes, tout en respectant les lignes directrices citées procure une réelle augmentation de performance. Les auteurs soulignent qu'un temps d'apprentissage est nécessaire, mais c'est un investissement rentable.

Les développeurs peuvent créer des interfaces utilisant le son, c'est un bon moyen de communication. L'interprétation de figures 3D peut être simplifiée par l'emploi d'environnements sonores structurés.

Les séquences sonores utilisées comme signalisation sont efficaces si elles sont construites judicieusement, en respectant certaines règles d'identification, de concordance conceptuelle, en déterminant judicieusement leurs paramètres physiques, choisies en fonction de préférences de l'utilisateur. Les informations cachées peuvent également être l'objet d'un traitement sonore particulier.

L'emploi du son au sein des interfaces homme-machine représente un véritable potentiel d'aide à la représentation d'objets complexes, comme les objets symboliques.

Comment le son peut-il aider l'utilisateur dans sa tâche d'analyse, d'interprétation, de compréhension des étoiles zoom ?

Comme nous l'avons fait dans le corps du chapitre, nous distinguons les sons utilisés comme moyen de signalisation de ceux utilisés comme moyen de perception, de communication d'informations.

### ***Le son : moyen de signalisation.***

L'idée d'un son typique pré-programmé est intéressante pour la mise automatique en évidence de certains événements.

*Les événements de l'environnement de travail :*

<b>Événement</b>	<b>Proposition de séquence sonore</b>
Ouverture d'une fenêtre de représentation	Ouverture d'une fenêtre (bruit de trafic)
Création d'une représentation étoile zoom	Démarrage d'une voiture
Suppression d'une représentation étoile zoom	Une feuille de papier chiffonnée
Fermeture de la fenêtre de représentation	Bris de vitre
Chargement d'un document enregistré	Chargement d'une arme à feu
Sauvegarde d'un document	Fermeture à clé d'une porte
Impression d'un document	Bruit de machine à écrire, d'imprimante
Quitter l'application	Fermeture violente d'une porte
Demande de l'aide en ligne	Exclamations d'une foule

Les séquences sonores sont associées aux icônes visuelles du logiciel. A l'image de ce qui existe pour certains logiciels sous WINDOWS (sur-impression d'un texte explicatif lorsque



la souris passe au-dessus de l'icône et s'y stabilise), ces séquences sonores sont audibles lorsque le curseur de la souris est situé au même emplacement.

Ces *earcons* font référence à des bruits familiers, le but est d'associer une action, un événement à un « environnement » sonore connu de l'utilisateur. D'autres événements, d'autres actions peuvent faire l'objet d'adjonctions sonores. L'utilisateur doit pouvoir, au moins, supprimer la génération sonore de ces séquences et, éventuellement les modifier ou en sélectionner d'autres en remplacement.

*Les événements propres aux étoiles zoom.*

Les types de variables

Associons un instrument à chaque type de variables :

Types de variables	Instruments associés
Quantitative sans intervalle	Trompette
Quantitative avec intervalle	Saxophone
Qualitative ordinale	Violon
Qualitative nominale	Guitare

Les instruments choisis sont aussi variés que possible. Leur identification en est plus aisée.

Les instruments associés aux variables quantitatives sont de la famille des instruments à vent tandis que les instruments à cordes sont liés aux variables qualitatives. Deux familles d'instruments pour deux ensembles de variables. Par la suite, nous respecterons cette typologie, chaque fois qu'une information concerne une variable qualitative ordinale, le violon est utilisé. Par défaut, une séquence sonore est associée à chaque axe. Il s'agit de la même séquence : une dizaine de notes d'une mélodie connue ou plus simplement les huit notes de la gamme : Do, Ré, Mi, Fa, Sol, La, Si, Do.

Les données manquantes

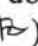
- Signaler les données manquantes

A la création d'une étoile zoom, signalons une donnée manquante par une séquence de notes. Pour représenter le nombre minimum de notes suggéré par SUMIKAWA (voir ci-avant), sept notes et une pause sont combinées. Les huit notes de la gamme sont jouées. Pour insister sur le fait qu'une donnée est manquante, la septième note est remplacée par une pause. Le rythme correspond à des croches. La séquence proposée remplace la séquence par défaut rattachée à l'axe.

La séquence émise est la suivante : Do, Ré, Mi, Fa, Sol, La, « pause », Do. La pause comme un trou, un manque dans la gamme pour souligner le fait qu'une donnée est absente pour une des variables représentées. A la création de l'étoile, la séquence sonore est jouée une seule fois. Elle est répétée si l'utilisateur clique sur l'axe de la variable.




- Déterminer le type de données manquantes

Déterminons le type de la données manquante en utilisant l'instrument approprié. Celui-ci identifiera le type de variables. Si pour une même représentation plusieurs données sont absentes, une icône graphique (par exemple un drapeau flottant dans le vent : ) indique la variable correspondante à la séquence jouée.

C'est par la couleur que nous identifions les données manquantes de type non-applicable. Aucune séquence sonore n'est utilisée pour les identifier des autres. Nous l'avons signaler plus haut, l'information sonore est éphémère, nous préférons signaler de façon stable et définitive le caractère non applicable de la variable pour un individu. Pour ne pas surcharger le son n'est pas utilisé.

### Les données particulières, aberrantes

Une icône visuelle dynamique (une petite bombe qui explose : ) balaye l'ensemble des axes en signalant visuellement et de manière sonore (le crépitement de la mèche qui brûle puis l'explosion), la présence d'une donnée particulière pour une variable de l'étoile représentée. Le balayage se fait automatiquement à la création de l'étoile ou par la sélection d'un item de menu.

Nous appelons valeur aberrante, particulière, une donnée dénotant par rapport à l'ensemble des données de la population sur la même variable. Par exemple : une variable qualitative nominale qui prenant deux valeurs simultanément alors que pour la même variable, la majorité des valeurs sont uniques; un intervalle particulièrement grand ou petit; une valeur extrême pour une variable, ...

Des petites analyses et sélections sont nécessaires avant de pouvoir activer cette fonctionnalité.


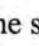
### La présence de variables particulières

Nous n'attachons pas de son particulier aux variables taxonomiques et aux dépendances. Une information visuelle existe déjà pour les identifier et les caractériser.

### La comparaison de deux étoiles zoom

Cette fonctionnalité est activable via un item de menu.

A la création de la seconde étoile zoom, deux curseurs (visuels) et l'émission de séquences sonores signalent à l'utilisateur les différences ou les similitudes importantes.

- Aux similitudes, nous associons l'icône visuelle  et une séquence sonore courte d'applaudissements.
- Aux différences, nous associons l'icône visuelle  et une séquence sonore courte d'huées.

Nous imaginons une icône visuelle (présente sur les deux représentations simultanément) qui balaye l'ensemble des axes en signalant, le cas échéant, visuellement et de manière sonore le type de différences ou de similitudes observables entre les deux étoiles.



Les comparaisons peuvent être marquées pour les situations suivantes :

- une donnée manquante (sur une étoile et pas l'autre);
- une différence important de la forme des intervalle d'une même variable;
- une variable à deux catégories et d'un individu à l'autre, les catégories sont différntes (exemple : la variable sexe)
- une différence importante de rang pour une variable ordonnée;
- ...

### *Le son : moyen de communication d'informations*

#### L'interprétation de la forme de l'étoile

Guidons l'utilisateur par le son dans sa localisation sur la représentation de l'étoile. Les déplacements de l'utilisateur selon l'axe horizontal sont détectés par une variation de la balance et ses mouvements verticaux par l'intensité d'émission.

Pour une étoile donnée, la balance se déplace horizontalement de gauche à droite; verticalement, l'intensité d'émission progresse de 10 à 20 décibels (voir figure 5.14.).

Cette localisation sonore est indépendante de la représentation de l'étoile, nous ne tenons pas compte de sa forme. Le but est de localiser le curseur de la souris de l'utilisateur. Ce dernier visualise l'étoile et le son l'aide à se repérer.

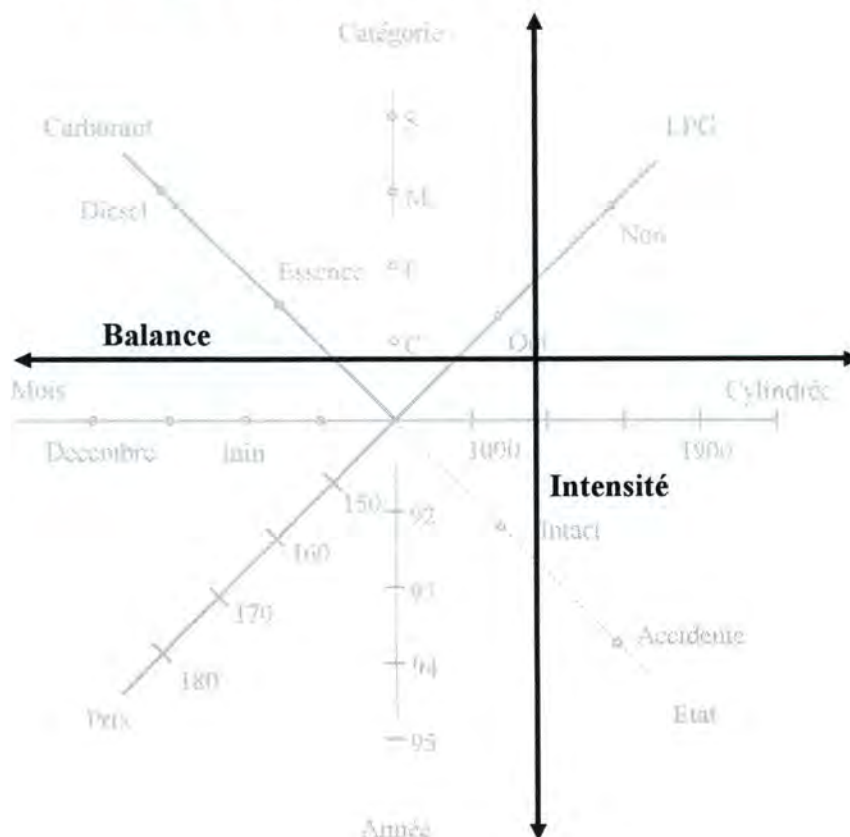


Figure 5.14. : utilisation de l'intensité sonore et de la balance comme aides à la localisation sur le graphique.



Nous suggérons l'emploi d'une séquence musicale continue de deux minutes. Les paramètres d'émission sont fonction de la position du curseur. L'émission terminée, l'utilisateur peut susciter une ré-émission de la séquence.

### L'interprétation de variables pondérées

Les variables pour lesquelles une pondération est utilisée sont représentées partiellement. Un axe en pointillé signale visuellement la caractéristique mais la représentation de l'étoile ne tient compte que de la valeur munie du poids le plus important. Le son pallie à cette carence.

#### • La représentation 2D

Nous proposons l'émission continue d'une note (le La (44 Hz)). En longeant l'axe de la variable, le rythme varie en fonction du poids. A l'image d'un monitoring de salle d'opération, le rythme s'accélère lorsque le poids augmente. D'un poids 0 à un, graduellement le rythme varie de 0 à 10 notes par seconde.

#### • La représentation 3D

Aidons l'utilisateur dans son interprétation de la représentation de la 3D des poids des variables. Décomposons en tranches égales les valeurs de poids potentielles souvent exprimées en pour-cent, créons cinq tranches égales de vingt pour-cent. Nous associons à chaque tranche un instrument différent (voir figure 5.15.).

La même mélodie est jouée en continu, un changement de tranche provoque un changement de l'instrument émetteur. Si le curseur se trouve à l'extérieur d'un bâtonnet de l'histogramme, l'émission est interrompue.

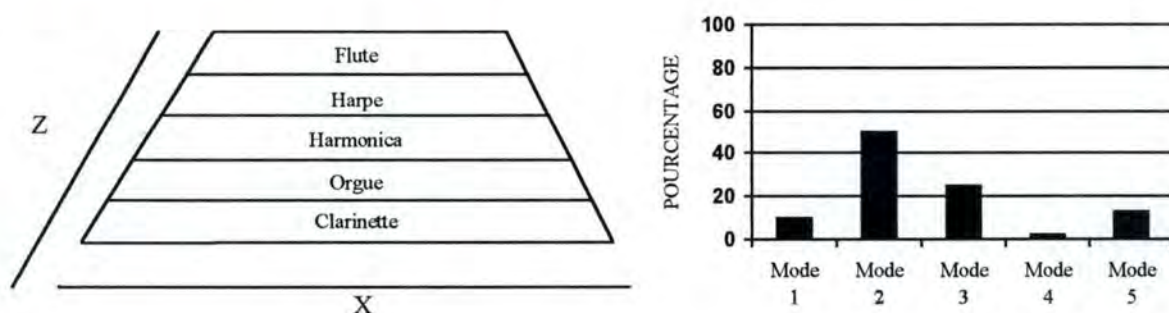


Figure 5.15. : répartition instrumentale en fonction de la valeur de poids.

Nous avons montré comment le son peut aider l'utilisateur dans sa tâche d'analyse, d'interprétation, de compréhension des étoiles zoom.

D'autres événements, actions ou phénomènes peuvent faire l'objet d'adjonctions sonores. L'utilisateur doit maîtriser l'émission sonore : il doit pouvoir au moins la supprimer, et éventuellement la modifier en en d'autres motifs de substitution.



### 1. Introduction

L'expérience de TAPP dans [TAPP 93] est menée dans le but de déterminer si la couleur et la taille de la police de caractères sont utiles dans l'affichage de code dans les tâches de programmation. Si c'est le cas, lequel de ces deux moyens est le plus rentable. L'hypothèse nulle de l'expérience est que ni l'un ni l'autre n'est bénéfique. Au terme de l'expérience, l'hypothèse nulle est réfutée. Il apparaît que de la couleur et de la taille de la police de caractères, le premier facteur est le préféré des utilisateurs. Nous nous limitons au détail des résultats concernant la couleur.

Voici en quatre étapes, la procédure d'expérimentation :

- 1) orientation des sujets, explication de la tâche à accomplir : dans un premier temps, réaliser un programme de test; par la suite, optimiser un autre programme volontairement modifié et erroné;
- 2) certains sujets (tirés au sort) exécutent la tâche avec un mécanisme d'affichage n'utilisant pas les couleurs, d'autres, issus d'un second groupe, utilisent un éditeur de texte couleur;
- 3) réalisation de la seconde tâche de programmation, phase d'optimisation; de nouveau, deux groupes de sujets sont créés;
- 4) les sujets complètent un questionnaire.

Au cours de l'expérience, plusieurs mesures sont relevées : temps de réalisation, nombre de tests (itération du programme),...

Résultats de la première phase (phase de programmation) :

	Couleur	Sans couleur
<b>Temps (Min.)</b>	14,167	16
<b>Itération</b>	5,5	2,769

Les différences sont statistiquement significatives. L'emploi de la couleur pour des tâches de programmation, améliore les performances.

Résultats de la seconde phase (phase d'optimisation et de correction) :

	Couleur	Sans couleur
<b>Temps (Min.)</b>	41,083	54
<b>Itération</b>	6,33	11

Ici aussi, les différences sont statistiquement significatives. L'emploi de la couleur pour des tâches de d'optimisation, de correction de programmes, améliore les performances.



L'utilisation de la couleur dans les tâches de programmation permet d'améliorer les performances. Les couleurs sont utilisées ici en mettant en évidence les commandes, et en variant la couleur d'affichage en fonction du contexte, du sujet (par exemple : le code concernant le tri est d'une couleur, les boucles d'une autre, les commandes d'une troisième).

Les réponses aux questionnaires sont unanimes pour l'usage de la couleur : les utilisateurs préfèrent travailler dans un environnement polychromatique.

Nous étendons ces conclusions aux tâches d'analyse. Depuis longtemps les couleurs sont utilisées au sein des applications interactives. Le matériel informatique actuel permet l'emploi de couleurs dans toutes les tâches habituelles. Les possibilités offertes par des périphériques de moins en moins onéreux (scanners, imprimantes couleurs, ...), autorisent la circulation de documents couleurs et un échange direct avec l'ordinateur. A l'inverse du son, l'emploi de la couleur est largement répandu dans toutes les applications et tous les domaines.

En plus du côté esthétique, l'apport de la couleur dans une application est significatif. Nous le constatons au travers de l'expérience de TAPP. Nombre d'autres expériences abondent dans ce sens.

Dans ce chapitre, nous avons montré l'utilité de la couleur par un exemple expérimental. Nous nous intéressons au « mécanisme » oculaire de perception des couleurs et à certains phénomènes d'illusion qu'il convient d'éviter. En fonction des trois paramètres que sont la texture, l'intensité et les nuances nous mettons en évidence des recommandations d'utilisation de la couleur.

Pour rédiger ce chapitre, nous nous sommes inspirés de : [VANDERDONCKT 94], [UNRUH 96], [TILO 96], [ROSE 96], [SHERRY 95], [ENKLAAR 95] et de [MACCHI 96].

## **2. La perception des couleurs**

### *Comment percevons-nous les différentes couleurs ?*

Deux faits :

- La couleur est trichromatique. Cela signifie que chaque teinte est une combinaison d'au moins trois longueurs d'onde. Leur combinaison est stable : le jaune composé de rouge et de vert se comporte exactement de la même façon s'il s'agissait d'une teinte pure. Si l'on combine une couleur supplémentaire à un assemblage existant, cela se passe comme si l'on prenait la couleur pure et que l'on opérât le mélange.
- Des effets de fatigue et de contraste sont facilement prévisibles. Lorsque l'on regarde longtemps une lumière rouge et que l'on déplace son regard sur une surface grise, une image résiduelle verdâtre gêne la vue. Un fait similaire se produit lorsque l'on fixe une lumière bleue, ici c'est une image résiduelle jaune qui se forme.



Beaucoup de théories tentent d'expliquer le fonctionnement de la vue couleur. La théorie de YOUNG-HELMHOLTZ et celle de HERING permettent ensemble d'expliquer de nombreux phénomènes de la perception des couleurs, mais ces théories ne peuvent les démontrer tous.

## **2.1. La théorie de YOUNG-HELMHOLTZ**

Sur la rétine de l'oeil, trois types de récepteurs répondent chacun le mieux à un certain intervalle du spectre visible. Ces récepteurs sont différents pour la lumière rouge, pour la lumière verte et pour la lumière bleue. La couleur d'une lumière provoque un échauffement des variable des trois récepteurs.

Une des faiblesses de cette théorie est qu'elle n'explique pas certains effets : par exemple la stabilité du jaune, l'image résiduelle.

RUSHTON en 1964, dans [ROSE 96], a montré que trois types de pigments dans les cellules de la rétine absorbent respectivement les longueurs d'ondes de 420 (bleu), 530 (vert) et 560 nm (rouge).

## **2.2. La théorie de HERING**

Dans le nerf optique situé entre la rétine et le cerveau, trois processus travaillent dans deux directions opposées : les directions anabolique et catabolique. Un processus accorde le rouge et le vert, un autre le jaune et le bleu et le dernier le noir et le blanc.

Rouge, jaune et blanc « empruntent » la direction anabolique, tandis que vert, bleu et noir s'associent à la direction catabolique.

Cette théorie permet d'expliquer pourquoi on ne peut percevoir un point bleu et jaune au même moment à la même place.

SVAETICHIN en 1956 dans [ROSE 96] a mis en évidence un potentiel électrique dans les cellules de la rétine d'un poisson. Potentiel positif lorsque les longueurs d'ondes courtes sont perçues et négatif pour les plus longues longueurs d'ondes.

DEVALLOIS en 1960 dans [ROSE 96] a découvert que certaines cellules, situées entre la rétine et le cortex du singe, s'échauffent plus ou moins en fonction des zones du spectre : des cellules pour le bleu et le jaune et des cellules pour le rouge et le vert.

Les deux théories prises en même temps permettent de dire que dans les cellules cônes de la rétine, trois types de pigments correspondent aux trois longueurs d'ondes 420, 530 et 560 nm. Le pigment bleu inhibe le processus bleu-jaune et excite l'autre, tandis que le pigment vert inhibe le processus vert-rouge et excite l'autre, le pigment rouge excite les deux processus.

Environ 10% des personnes souffrent d'une perception déficiente ou anormale des couleurs. Ainsi le protonope ne distingue pas le rouge, le deuteronope pas le vert, le tritanope pas le bleu. L'anomale est atteint d'une diminution de la sensibilité à la couleur (une ou plusieurs), l'acromate ne perçoit pas les couleurs (faute de cônes).



### 3. Les illusions

Le système visuel humain est sensible aux illusions, celles-ci affectent la perception. Il est opportun de rappeler certains effets pour éviter de créer des représentations mal interprétables.

Une image n'est pas perçue point par point mais en termes de formes connectées. De la perception de ces formes géométriques, dépend l'interprétation que l'utilisateur en donnera.

Des contradictions visuelles renforcées et de complexité réduite améliorent le niveau de perception des formes géométriques.

Nous illustrons certains types d'illusions. Dans la mesure du possible, il faut l'éviter ou tout au moins en minimiser l'impact.

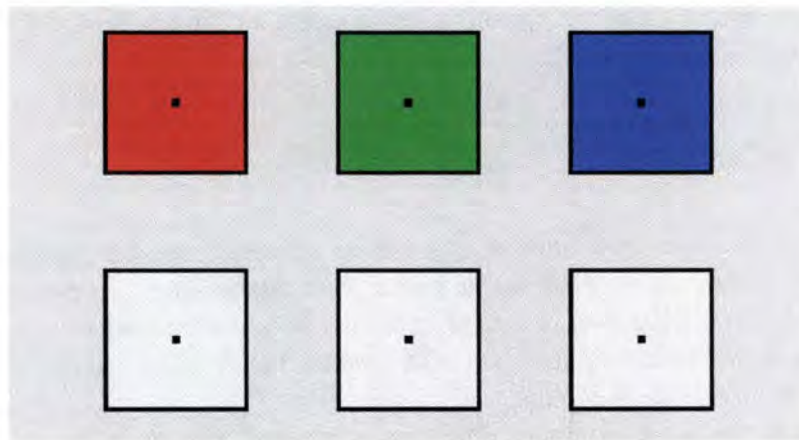


Figure 6.1. : Les images résiduelles.

Si l'on focalise son attention sur le carré supérieur (figure 6.1.) pendant 15 secondes, que l'on déplace son regard vers le carré inférieur correspondant, une image résiduelle se forme. L'image résiduelle correspond exactement aux couleurs complémentaires. Ce type d'illusion ne peut être éliminé.

La perception de l'intensité visuelle d'un pixel par rapport à son intensité absolue est fortement influencée par son voisinage dans l'image. Si des représentations visuelles doivent être interprétées quantitativement selon leur intensité, des erreurs apparaissent.



Figure 6.2. : Illusion d'intensité, les deux carrés rouge ont exactement la même intensité. Pourtant la carré rouge de droite apparaît plus foncé.



Deux mesures à prendre pour contrecarrer cet effet d'illusion :

- Uniformiser le fond : les comparaisons relatives seront toujours possibles. C'est ce que nous proposons pour les étoiles zoom : un fond blanc, le même pour toutes les représentations.

- Eviter de quantifier des valeur par des nuances d'une même couleur. Nous suggérons de ne pas employer cette technique de travail : les variations d'intensité au sein d'une même couleur ne sont pas préconisées.

Le phénomène d'irradiation est causé lorsque des objets clairs sont centrés sur un fond foncé. L'objet clair apparaît alors plus grand (figure 6.3. ci-dessous).

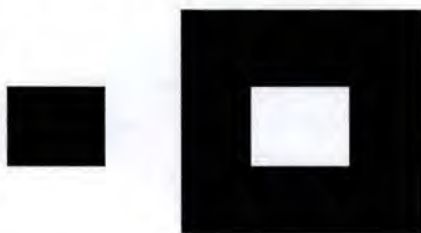


Figure 6.3. : Le phénomène d'irradiation.

Nous incitons l'emploi d'un fond blanc pour la représentation des étoiles zoom.

## 4. Les couleurs et les étoiles zoom

### 4.1. La texture



Les zones de l'écran où la couleur est la plus intense et la plus saturée captent immédiatement l'attention de l'utilisateur [VANDERDONCKT 94]. Pour la représentation des étoiles zoom, nous préconisons l'emploi des couleurs sans texture bigarrée, la pleine couleur.

### 4.2. L'intensité



ENKLAAR recommande l'utilisation de la variation de l'intensité monochromatique. Selon lui, la perception est plus rapide et donc plus aisée. Nous pensons que l'emploi simultané de plusieurs couleurs différentes facilite l'identification et la compréhension de la représentation. Les différents éléments d'une figure sont plus aisément, plus rapidement identifiables s'ils sont colorés différemment. Nous proposons de jouer sur les nuances dans les couleurs plutôt que sur les nuances d'une couleur. Nous suggérons de se servir de différentes couleurs et non du contraste d'une couleur lors de la représentation d'étoile zoom.



### 4.3. Les nuances



En général, le concepteur d'une interface utilise peu de couleurs différentes. Passons en revue les principales et déterminons-en les implications pour la représentation des étoiles zoom.

#### 4.3.1. Le bleu

J. VANDERDONCKT déconseille le bleu pour les textes, les traits fins et les petits objets. De plus, il est déconseillé de présenter simultanément des nuances dans le bleu. Le bleu est une couleur réservée aux données moins importantes. Nous avons choisi de ne pas utiliser le bleu dans les représentations de l'étoile zoom. D'autres couleurs plus efficaces suffisent.

#### 4.3.2. Le vert

Cette couleur est intéressante pour la mise en évidence. Elle signale en particulier un état normal pour un paramètre quelconque. Ce n'est pas dans ce sens que nous l'employons. Nous préconisons d'employer le vert pour représenter les données non-applicables<sup>1</sup>. Nous avons choisi d'utiliser le vert qui se marie bien avec le rouge des étoiles. Cette couleur nous permet la mise en évidence des variables non-applicable sans heurter l'utilisateur. Il ne s'agit pas d'une erreur, l'attention de l'utilisateur ne doit pas être captée violemment.

#### 4.3.3. Le jaune

Cette couleur est utilisée pour avertir. Nous ne l'employons pas.

#### 4.3.4. Le rouge

Nous proposons de colorer les étoiles en rouge. J. VANDERDONCKT préconise l'emploi du rouge pour alarmer, signaler les erreurs, signifier des arrêts. Puisque l'attention de l'utilisateur doit se focaliser sur la représentation de l'objet symbolique nous proposons cette couleur.

Notre suggestion est renforcée par le conseil suivant de J. VANDERDONCKT : pour les grosses épaisseurs, utiliser le noir, le bleu, le rouge, ... . Nous le verrons, le noir est réservé aux informations sur le contexte. Le bleu ne capte pas l'attention comme le rouge. C'est vers le rouge que s'est porté notre choix. Le contraste entre cette couleur et le fond d'écran est important : sur un fond clair (nous préconisons l'emploi du fond blanc), ce sont les couleurs : le noir, puis le bleu, puis le rouge qui ressortent le mieux. De nouveau notre proposition se justifie.

---

<sup>1</sup> Elles correspondent à des variables qui dans certaines circonstances ne peuvent avoir de valeurs. Une variable mesurant la longueur de la période de grossesse n'a pas de sens chez les hommes. Ce n'est donc pas une donnée manquante au sens d'une erreur de mesure, mais une impossibilité pour une variable de prendre une valeur compte tenu du contexte.



#### **4.3.5. Le noir**

Nous l'avons déjà évoqué, le noir est réservé à tous les éléments du contexte : les axes, les graduations, les textes, les libellés, les légendes, ... . Conventionnellement, cette couleur est utilisée par défaut. Classiquement la police de caractère est noire. En plus des « us et coutumes », l'emploi du noir se justifie pleinement pour la représentation des objets de contour d'épaisseur fine et/ou pour des objets peu épais. Le noir est couramment utilisé comme couleur de base.

Pour la représentation des axes des variables à données manquantes, nous introduisons une entorse à notre souhait : ne pas jouer sur les nuances d'une seule couleur. Nous préconisons l'emploi du gris clair pour localiser la donnée manquante. Nous suggérons de travailler avec 70-80% de gris. Ce pourcentage permet de ne pas dénoter dans l'ensemble de l'étoile (de ne pas gâcher la symétrie de l'étoile). Il est suffisamment sensible pour porter rapidement l'attention de l'utilisateur sur le phénomène.

#### **4.3.6. Le blanc**

C'est certainement ce qui est le plus répandu, le plus classique. De plus, le fond de teinte blanche déforme le moins la taille d'un objet en avant plan.

### **5. Conclusion**

La couleur est un moyen largement utilisé pour communiquer de l'information. Par sa couleur, un objet prend directement une signification (rouge pour le chaud, bleu pour le froid, ...). Le choix des couleurs de représentation est important. Pour identifier les divers éléments d'une figure, les couleurs différentes permettent de gagner en efficacité, en simplicité. Plus la densité d'un dessin est grande, plus les différentes couleurs jouent un rôle prépondérant dans l'identification des données (J. VANDERDONCKT). Le temps de recherche d'une donnée bien colorée par rapport à ses voisines est considérablement diminué si le codage est subtil.

L'emploi des couleurs accroît l'attrait, la crédibilité, la compréhensibilité des données à mettre en évidence dans un contexte particulier. Actuellement, le matériel disponible permet de respecter au mieux toutes ces recommandations. Tous les utilisateurs cibles de la représentation en étoile zoom sont équipés d'écrans couleur. Le prototype conçu par M. NOIRHOMME-FRAITURE et M. ROUARD s'adresse à des utilisateurs familiers de l'outil informatique, il est opportun de respecter les conventions admises de tous.

Comme pour le son, il est important que l'utilisateur reste maître de son environnement de travail. Il doit pouvoir modifier des paramètres. Le concepteur permettra à l'utilisateur de s'accorder à ses propres préférences en matière de couleurs soit de laisser la possibilité de modifier tous les composants de l'environnement d'affichage.



### 1. Introduction

M. NOIRHOMME et M. ROUARD dans [NOIRHOMME-FRAITURE 97b] soulignent l'intérêt de la combinaison 2D et 3D : chacune d'elles a des avantages et les informations qu'elle permet de représenter intéressent l'utilisateur. Si les objets 3D sont plus difficiles à interpréter (la création d'image mentale n'est pas simple), ils permettent d'afficher plus d'informations. La difficulté n'est pas seulement liée au passage à une troisième dimension mais également à l'apport informatif supplémentaire.

Pour palier à cette difficulté, comme nous l'avons vu au chapitre 05, le son aide à interpréter des blocs diagrammes 3D : MEREU *et al.* dans [MEREU 96] emploient le son pour fournir aux utilisateurs des moyens supplémentaires de compréhension.

Nous insistons dans ce chapitre sur les raisons de l'utilisation de la 3D. Lors de la représentation 3D, certains éléments se retrouvent cachés par d'autres. Nous abordons ces situations de visibilité réduite en proposant des voies de contournements. Des outils particuliers permettant la manipulation des représentations 3D sont analysés.

### 2. Les circonstances d'emploi de la 3D

#### 2.1. Introduction

La représentation graphique de données quantitatives est de plus en plus répandue dans tous les domaines : sciences, économie et même dans la presse. Puisque la technologie (*hardware* et *software*) le permet, cette représentation graphique devient chaque jour plus accessible et variée. Il y a peu, l'affichage 3D se limitait à quelques domaines particuliers : CAD, analyses scientifiques (représentation de molécules). Plus récente et plus attrayante, la représentation 3D est employée comme moyen d'illustration. Soulignons l'intérêt réel de la représentation 3D, la justification de son emploi.

#### 2.2. Les expériences

LEVY *et al.* dans [LEVY 96] ont voulu déterminer les préférences 2D-3D des utilisateurs.

Contexte de l'expérience : les sujets doivent s'identifier à un chercheur scientifique qui doit présenter le résultat de ses recherches à un comité de direction. Comment préparer les canevas de graphiques (2D et 3D) qu'il compte présenter.



Les objectifs de l'exercice sont les suivants :

- présenter les points importants et la structure de l'étude;
- ne montrer que certains points particuliers;
- se préparer à répondre à des questions de détail à l'aide des graphiques;
- montrer les tendances et les extrapolations;
- identifier certaines phases de l'étude limitées à un ou deux points;
- faciliter aux interlocuteurs la mémorisation des points essentiels.

Dans le cadre de leurs recherches, LEVY *et al.* mènent deux autres expériences de confection de graphique et de comparaison 2D - 3D.

En matière de graphiques 3D, voici les préférences des sujets interrogés :

- les histogrammes 3D sont les plus plébiscitées des représentations tridimensionnelles;
- les vues du dessus sont largement préférées à celles du dessous;
- les graphiques linéaires 3D : la représentation en volumes est préférée à celles en ligne;
- les graphiques 3D sont plus utiles pour identifier certains points de détail : par contre, ils expriment moins bien les tendances;
- les graphiques 3D sont plus efficaces pour la mémorisation des informations;
- les histogrammes 3D permettent mieux les comparaisons.

Pour leurs « utilisations personnelles », les sujets préfèrent l'emploi de graphiques 2D par contre pour la transmission de résultats à d'autres personnes, ils préconisent l'utilisation de la 3D.

### 2.3. Synthèse

De ces expériences et des constatations déjà mentionnées, nous proposons l'utilisation de la 3D en plus ou à la place de la 2D pour les raisons suivantes :

- mise en évidence des **détails**;
- **transmission** entre utilisateurs de l'information;
- favoriser la **mémorisation** de la représentation;
- privilégier les **vues plongeantes** (mentalement aisées à interpréter);
- intérêt de l'emploi d'**histogrammes** pour les comparaisons (graphiques en barres).

## 3. 3D et visibilité

L'intérêt de la 3D dans le domaine de l'architecture et de la construction est plus que justifié. Le facteur primordial de cet intérêt réside dans la concordance de trois mesures orthogonales d'un volume à représenter. Dans ce cas les trois mesures sont de même nature (des distances) et de même échelle. Cette constatation ne s'applique pas aux étoiles zoom. Une étoile zoom assure la représentation d'un certain nombre de variables qui bien que liées entre elles ne doivent pas nécessairement être représentées à la même échelle.



L'affichage à l'écran d'une représentation 3D génère *de facto* des zones cachées. Plusieurs techniques permettent de contourner ce problème de manque de visibilité :

- les distorsions de la représentation;
- les déplacements de la représentation.

COWPERTHWAITE *et al* dans [COWPERTHWAITE 96] décrivent une solution au problème de l'occlusion de zones dans la représentation d'objets 3D. Une fonction de distorsion est utilisée pour mettre en évidence des éléments cachés dans les visualisations en trois dimensions.

Dans la représentation 3D, il est courant qu'un objet soit partiellement ou totalement dissimulé par d'autres. La représentation 3D d'une molécule chimique complexe ne permet pas toujours de visualiser tous les atomes. Les auteurs proposent de modifier la géométrie 'centrale' de l'objet pour mieux visualiser le coeur. La méthode consiste schématiquement à choisir une ligne de vue, à calculer un déplacement des points voisins à cette ligne. La distorsion est plus ou moins forte selon le souhait de l'utilisateur. La déformation est d'autant plus marquée que l'on se trouve à proximité de la ligne de vue. La mise au point de l'algorithme permet qu'il puisse être appliqué à plusieurs éléments d'un même objet.

Si l'intérêt de la solution présentée par COWPERTHWAITE *et al* est évident pour la représentation de molécules chimiques, il est beaucoup moins net appliqué aux objets où la géométrie joue un rôle d'importance comme en architecture ou en charpente métallique ou pour les objets symboliques.

Néanmoins, son utilisation pour la mise en évidence de certains aspects cachés de l'objet symbolique reste utile. Le déplacement d'angle de vue permis par de nombreux logiciels n'autorise que la variation de la vision d'un objet depuis l'extérieur. Dans les meilleurs cas l'accès virtuel au sein des objets est possible. Les auteurs proposent de visualiser le coeur des objets en appliquant une fonction de distorsion, tout en restant 'hors' de celui-ci.

Le déplacement des représentations permet l'apparition des éléments dissimulés. La rotation de la figure autour d'un axe central passant par l'origine de tous les axes fait apparaître les éléments de la partie arrière de l'étoile zoom. Jumelés avec la possibilité de modifier l'angle de vue d'une étoile, les déplacements offrent l'opportunité d'analyser l'étoile zoom sous tous ses angles.

Lorsque l'utilisateur compare plusieurs étoiles zoom, les déplacements (rotations et modifications de l'angle de vue) sont simultanément appliqués à toutes les représentations. Pour que les comparaisons soient permises, cette fonctionnalité est obligatoire.

Les situations où des éléments de la représentation 3D sont dissimulés derrière d'autres représentent les limites de la représentation 3D. Néanmoins, ces difficultés sont contournées. La rotation et la modification de l'angle de vue sont retenus dans la cadre des représentations 3D des étoiles zoom.



## 4. Les manipulations

L'intérêt de la 3D dans le domaine architectural n'est plus à justifier. SAKAI dans [SAKAI 96] décrit un outil permettant à plusieurs utilisateurs la manipulation d'objets 3D : *Flying Finger*. Deux utilisateurs travaillent simultanément sur le même objet, ils utilisent le contrôle de l'angle de vue en se basant sur des coordonnées sphériques à partir d'une souris à déplacements linéaires. Les potentialités d'aide à la présentation publique de représentation 3D sont réelles : chacun manipulant l'objet à sa guise.

Les moyens classiques (souris et clavier) présents actuellement dans les systèmes d'aide à la conception ne permettent pas aux utilisateurs des manipulations directes et intuitives des formes 3D complexes.

Pour résoudre ce problème, MURAKAMI dans [MURAKAMI] propose un nouveau système d'interface. Il s'agit d'un objet élastique se présentant sous forme d'un volume : un cube. En déformant ce mécanisme, l'utilisateur manipule, déforme un volume 3D représenté à l'écran. L'auteur décrit un prototype confectionné en polyuréthane.

Il s'agit d'un cube de 10 cm de côté, chaque arête et diagonale constitue un capteur déformable. Les déformations sont transmises, via un port série, à un logiciel qui les interprète et les représente sur la forme affichée.

Ce type de manipulateur aide l'utilisateur dans son appréhension de l'objet 3D. Plutôt que de déplacer le curseur via la souris, l'utilisateur dispose d'une étoile à 16 branches pour ses manipulations. A l'image de ce que décrit MURAKAMI, l'outil en mousse est munis de capteurs. Une pression correspond à un clic souris et déclenche des événements et actions tels que ceux présentés dans le chapitre sur le son. Le déplacement du doigt le long d'une branche de l'étoile permet par exemple d'apprécier de manière auditive l'interprétation sonore des poids d'une variable pondérée. Si une étoile ne représente que quatre variables (4 axes), seuls les quatre bras correspondants de l'étoile sont actifs. L'utilisateur manipule cet outil comme un « joystick ».

Différents auteurs soulignent dans leurs publications l'emploi d'une souris spéciale permettant des déplacements dans les trois dimensions. Le temps d'adaptation à ce nouveau mode de manipulation est considérable. L'effort d'apprentissage est conséquent. Les résultats ne sont pas à la hauteur des attentes. Pour en avoir testé une et en accord avec les auteurs, nous suggérons de limiter, hormis notre proposition d'étoile 3D en mousse, les moyens d'interaction à la souris et au clavier.



## 5. Conclusion

L'apport de la troisième dimension se justifie pour les objectifs suivants :

- mise en évidence des **détails**;
- **transmission** entre utilisateurs de l'information;
- favoriser la **mémorisation** de la représentation.

Lorsque le concepteur intègre la 3D dans les représentations, il convient de :

- préférer les **vues plongeantes**;
- d'employer les **histogrammes** pour les comparaisons (graphiques en barres).

Les situations où des éléments de la représentation 3D sont dissimulés derrière d'autres représentent les limites de la représentation 3D. Dans la cadre des représentations 3D des étoiles zoom ces difficultés sont contournées : la rotation et la modification de l'angle de vue sont retenus.

L'utilisateur dispose d'une étoile à 16 branches pour ses manipulations. L'outil en mousse est munis de capteurs. Une pression correspond à un clic souris et déclenche des événements et actions tels que ceux présentés dans le chapitre sur le son. Le déplacement du doigt le long d'une branche de l'étoile permet par exemple d'apprécier de manière auditive l'interprétation sonore des poids d'une variable pondérée. L'utilisateur manipule cet outil comme un « joystick ».



## 1. Introduction

Nous avons vu au chapitre 03, qu'il est possible de représenter graphiquement les objets symboliques. Mais représenter des objets symboliques par  $p$  variables peut ne pas être réalisable si  $p$  est très grand. De plus, la représentation munie de  $p$  variables n'est peut-être pas judicieuse : des variables sont fortement corrélées et, dessiner toutes les variables encombre le graphique ne permettant pas une vision claire et analytique.

Dans ce chapitre, nous présentons des méthodes ou des critères de sélection des variables. Le but à atteindre est d'identifier les variables qui, ensemble, permettent de mieux discriminer les individus d'une population. Pour ce faire, nous nous sommes inspirés de [AGRESTI 90], [JOHNSON 92], [MOLENBERGHS 96a], [MOLENBERGHS 96b], [NETER 96] et de l'aide en ligne de certains logiciels.

Le principe général de la sélection de variables est d'estimer des critères ou paramètres, pour un modèle regroupant une ou plusieurs variables et de modifier ce nombre de variables. Par exemple, on commence avec une variable, puis progressivement on en ajoute d'autres, le choix de celles-ci étant guidé par les critères fixés. Parfois, et c'est le cas pour les critères que nous présentons ci-dessous, on les calcule pour toutes les combinaisons de variables possibles. En fonction de ces critères on détermine la meilleure combinaison, le meilleur modèle.

Dans le cas de la représentation graphique d'objets symboliques, en retenant à 16 comme nombre maximum d'axes représentables sur une étoile zoom, on se choisira le meilleur modèle contenant au plus 16 variables en fonction d'un critère ou de la méthode. Pour la représentation d'objets symboliques d'un même type (par exemple : les voitures d'occasion), ce sont toujours les mêmes variables qui seront représentées et dans le même ordre sur l'étoile.

Les méthodes que nous présentons ci-dessous, s'adressent en particulier à des variables quantitatives. Les variables des objets symboliques ne sont pas nécessairement de ce type. Des opérations de codage et de transformation s'imposent pour permettre l'application des méthodes décrites. Nous reprenons à l'annexe 1 les opportunités de modifications. Ces remarques sont valables pour les chapitres 08 et 09 traitant respectivement de la sélection des variables et des données manquantes.

## 2. Des critères de sélection

Différents critères de comparaison de modèles peuvent être utilisés, présentons les suivants :  $R_p^2$ ,  $MSE_p$ ,  $C_p$  et  $PRESS_p$ .

### 2.1. Le critère $R_p^2$ ou $SSE_p$

Le critère  $R_p^2$  est en fait une version du coefficient  $R^2$  de détermination multiple, utilisé pour identifier plusieurs « bons » sous-ensembles de variables. Plus ce critère est élevé, meilleur sera le modèle (voir tableau 8.1.).  $R_p^2$  indique qu'il y a  $p$  paramètres ou  $p-1$  variables présentes dans la fonction de régression pour laquelle  $R_p^2$  est calculé.



L'emploi de  $R_p^2$  est équivalent à l'utilisation de  $SSE_p$ , la somme des carrés des erreurs. Avec  $SSE_p$ , plus le critère est petit, meilleur sera le modèle (voir tableau 8.1.).

La relation entre  $R_p^2$  et  $SSE_p$  est la suivante :

$$R_p^2 = 1 - \frac{SSE_p}{SSTO}$$

Lorsque l'on utilise le critère  $R_p^2$ , on ne cherche pas à obtenir la plus grande valeur possible, on calcule la différence entre  $R_p^2$  et  $R_{p+1}^2$ . Tant que cette valeur n'est pas minime, n'atteint pas un seuil fixé, on continue la recherche de modèles par introductions successives de nouvelles variables.

Remarque :

$$SSTO = SSE + SSR$$

Ou encore :

$$Total\ sum\ of\ squares = Error\ sum\ of\ squares + Regression\ sum\ of\ squares$$

Ou encore :

$$\sum_i (Y_i - \bar{Y})^2 = \sum_i (Y_i - \hat{Y}_i)^2 + \sum_i (\hat{Y}_i - \bar{Y})^2$$

## 2.2. Le critère $MSE_p$ ou $R_a^2$

Puisque  $R_p^2$  ne tient pas compte du nombre de paramètres dans le modèle de régression, l'emploi d'un coefficient ajusté de détermination multiple  $R_a^2$  est suggéré comme alternative :

$$R_a^2 = 1 - \left( \frac{n-1}{n-p} \right) \frac{SSE}{SSTO}$$

Le critère tient compte du nombre de paramètres à estimer dans le modèle de régression, et ce au travers du nombre de degrés de liberté. Les utilisateurs de  $MSE_p$  cherchent à identifier quelques modèles pour lesquels le critère  $MSE_p$  est minimum ou proche du minimum tel que : ajouter plus de variables ne vaut pas la peine. Dans ce cas aussi, c'est la différence avec le critère de l'étape précédente qui est calculée et qui sert d'élément de décision.



### 2.3. Le critère $C_p$

Le modèle qui inclut toutes les  $P-1$  variables potentielles est supposé être choisi pour que  $MSE(X_1, \dots, X_{P-1})$  soit un estimateur le moins biaisé possible de  $\sigma^2$ .

$MSE$  est la moyenne des carrés des erreurs, écarts entre les prédictions et la moyenne générale.

Le critère  $C_p$  se calcule de la façon suivante :

$$C_p = \frac{SSE_p}{MSE(X_1, \dots, X_{p-1})} - (n - 2p)$$

Lorsqu'on utilise le critère  $C_p$ , on cherche à identifier un sous-ensemble de variables pour lequel :

- la valeur de  $C_p$  est petite et,
- la valeur de  $C_p$  est proche de  $p$ .

Dans ce cas, le biais du modèle de régression est petit (tableau 8.1.).

### 2.4. Le critère $PRESS_p$ .

Le critère  $PRESS_p$  (prédiction de la somme des carrés des écarts) est une mesure de la qualité de prédiction d'un modèle.

$$PRESS_p = \sum_{i=1}^n (Y_i - \hat{Y}_{i(i)})^2$$

Où :

$\hat{Y}_{i(i)}$  est la valeur prédite, pour l'individu  $i$ , par le modèle de régression construit en supprimant le  $i^{\text{ème}}$  élément du set de données de départ, c'est-à-dire sur les  $n - 1$  autres individus.

$PRESS_p$  diffère généralement peu de  $SSE_p$ , particulièrement dans le cas où le nombre d'individus dans le modèle est grand ou si le nombre de variables est grand.

Les modèles avec de petites valeurs de  $PRESS_p$  sont considérés comme de bons modèles candidats (voir tableau 8.1.).

Afin d'illustrer ces propos théoriques, nous présentons au tableau 8.1. le résultat du calcul de ces différents critères pour les modèles possibles dans le cas d'une étude médicale portant sur 54 individus caractérisés par 4 variables.



Variables dans le modèle	$p$	$df$	$SSE_p$	$R_p^2$	$MSE_p$	$C_p$	$PRESS_p$
Aucune	1	53	3,9728	0	0,0750	1721,6	4,1241
$X_1$	2	52	3,4961	0,120	0,0672	1510,8	3,8084
$X_2$	2	52	2,5763	0,352	0,0495	1100,1	2,8627
$X_3$	2	52	2,2153	0,442	0,0426	939,0	2,4268
$X_4$	2	52	1,8776	0,527	0,0361	788,2	2,0292
$X_1, X_2$	3	51	2,2325	0,438	0,0438	948,7	2,6388
$X_1, X_3$	3	51	1,4072	0,646	0,0276	580,2	1,6095
$X_1, X_4$	3	51	1,8758	0,528	0,0368	789,4	2,1203
$X_2, X_3$	3	51	0,7430	0,813	0,0146	283,7	0,8352
$X_2, X_4$	3	51	1,3922	0,650	0,0273	573,5	1,5833
$X_3, X_4$	3	51	1,2453	0,687	0,0244	507,9	1,4287
$X_1, X_2, X_3$	4	50	0,1099	0,972	0,0022	3,1	0,1405
$X_1, X_2, X_4$	4	50	1,3905	0,650	0,0278	574,8	1,6513
$X_1, X_3, X_4$	4	50	1,1156	0,719	0,0223	452,0	1,3286
$X_2, X_3, X_4$	4	50	0,4652	0,883	0,00930	161,7	0,5487
$X_1, X_2, X_3, X_4$	5	49	0,1098	0,972	0,00224	5,0	0,1456

Tableau 8.1. : résultats des calculs des différents critères pour les modèles possibles dans le cas d'une étude médicale portant sur 54 individus caractérisés par 4 variables [NETER 96].

#### Remarque :

pour que les calculs des critères et méthodes soient réalisables, il faut que  $n < p$ .

### 3. Les procédures automatiques de sélection de variables

#### 3.1. Introduction

Des algorithmes existent pour calculer, déterminer, en fonction de l'un ou l'autre critère (voir ci-dessus), les meilleurs des sous-ensembles de variables. Mais lorsque ce nombre de variables est fort élevé (plus de 30 - 40), utiliser ces algorithmes ne paraît plus réalisable. Une procédure automatique de sélection des variables est pourtant fort utile. Les procédures *Forward Stepwise* et *Backward Stepwise* sont certainement les plus répandues.

En général, la méthode est la suivante : développer une séquence de modèles de régression en ajoutant ou en retirant à chaque étape une variable, selon un critère et un seuil définis par l'utilisateur.

Le critère d'ajout ou de retrait de variables peut être la réduction de la somme des carrés des erreurs, le coefficient de corrélation partiel,  $t^*$  ou  $F^*$ .

La différence principale entre les méthodes utilisant les critères vus ci-avant et les procédures automatiques est que ces dernières se terminent généralement avec l'identification d'un modèle de régression unique, modèle considéré comme le meilleur. Les deux méthodes sont parfois combinées afin d'optimiser le résultat.



### 3.2. Forward Stepwise Regression

Nous décrivons la procédure *Forward Stepwise*, son algorithme utilisant  $F^*$ .

1. La routine *Forward Stepwise* estime d'abord un modèle de régression linéaire pour chaque variable de l'étude. Pour chaque modèle  $F^*$  est calculé :

$$F_k^* = \frac{MSR(X_k)}{MSE(X_k)}$$

Où :

$MSR$  et  $MSE$  expriment respectivement la moyenne des carrés des erreurs résiduelles et la moyenne des carrés des erreurs, c'est à dire  $\frac{SSR}{p-1}$  et  $\frac{SSE}{n-p}$ , avec  $p-1$  (le nombre de variables dans le modèle) et  $n-p$  les nombres de degrés de liberté. A cette étape (1),  $p=2$ .

La variable avec la plus haute valeur de  $F^*$  est candidate pour la première sélection.

2. Supposons que la variable  $X_i$  est sélectionnée à l'étape précédente, la routine va maintenant estimer tous les modèles avec deux variables, où  $X_i$  est membre de la paire. Pour chaque modèle calculé,  $F_k^*$  est déterminé :

$$F_k^* = \frac{MSR(X_k | X_i)}{MSE(X_i, X_k)}$$

La variable avec le  $F_k^*$  le plus élevé est sélectionnée. Si la valeur de  $F^*$  dépasse un seuil fixé, la seconde variable est ajoutée et le programme se termine.

3. Si l'on suppose que la variable  $X_j$  est ajoutée à la deuxième étape, la procédure *Forward Stepwise* examine alors si une des variables présentes dans le modèle peut être abandonnée. On calcule :

$$F_i^* = \frac{MSR(X_i | X_j)}{MSE(X_j, X_i)}$$

Si le plus petit  $F_i^*$  est inférieur à un seuil prédéterminé alors la variable peut être abandonnée.

4. Supposons que  $X_i$  et  $X_j$  soient dans le modèle après l'étape 3, la procédure calcule maintenant laquelle des variables restantes est la prochaine candidate à la sélection. Les deux dernières étapes sont répétées.

#### Terminaison :

La routine s'arrête lorsque le seuil  $F_k^*$  est dépassé la première fois lors de l'étape d'addition d'une variable.



**Remarque :**

la routine permet qu'une variable ajoutée au modèle lors d'une étape, soit supprimée lors d'une des étapes suivantes si son apport au modèle, n'est plus jugé significatif.

### **3.3. D'autres procédures automatiques de recherche**

D'autres procédures automatiques de recherche sont utilisables pour déterminer un bon sous-ensemble de variables utiles. Nous en mentionnons deux.

#### **3.3.1. Forward Selection**

C'est une version simplifiée de la régression *Forward Stepwise*; ici on omet de tester si une variable déjà entrée dans un modèle peut en être éjectée.

#### **3.3.2. Backward Selection**

C'est une méthode opposée à la *Forward Stepwise*. Elle commence avec le modèle contenant toutes les variables et identifie celle qui a la plus petite valeur de  $F^*$ . Si cette valeur est inférieure à un seuil fixé, la variable est abandonnée. Une procédure *Backward Elimination* peut aussi être adaptée pour que une variable éliminée plutôt puisse être récupérée ultérieurement, cette modification est appelée la procédure de régression *Backward Stepwise*.

## **4. Conclusions**

Nous nous sommes intéressés à une série de critères et méthodes qui permettent de choisir significativement les variables à conserver. Le but est de permettre la représentation graphique d'objets symboliques. Cette représentation, pour être claire et efficace, ne doit pas contenir trop d'axes, nous pensons à un maximum de 16. Il s'agit, lorsque le nombre de variables, des données de base, est trop élevé ( $> 16$ ) de sélectionner celles qui permettent de discriminer au mieux les individus. Différentes méthodes efficaces existent et sont couramment implémentées et rodées.

Ces méthodes constituent des outils d'aides complémentaires; pratiquement, l'utilisateur sélectionne les variables qu'il désire analyser et représenter graphiquement.



### 1. Introduction

Lors d'enquêtes, des individus interrogés ne répondent pas ou de façon partielle voire incohérente. Lors de prises automatiques de mesures continues, un appareil défectueux peut biaiser certaines données.

Ces données manquantes ne doivent pas altérer les traitements des analyses; il faut soit les éliminer, soit les remplacer.

Dans ce chapitre, nous nous intéressons à différentes méthodes de résolution de ce genre de situations. Pour ce faire, nous nous sommes inspirés de [AGRESTI 90], [JOHNSON 92], [MOLENBERGHS 96a], [MOLENBERGHS 96b], [NETER 96] et de l'aide en ligne de certains logiciels.

Ce type de manipulations sur les données a des implications sur les interprétations et conclusions issues des analyses. Honnêteté scientifique et rigueur s'imposent.

En toute généralité, la présence de données manquantes implique directement une augmentation de l'imprécision dans l'analyse et dans les résultats, imprécision dépendant du nombre de données manquantes. Il y a lieu d'en tenir compte dans l'interprétation des résultats.

### 2. Les données manquantes

Souvent, des observations du set de données sont manquantes. Les raisons de cette carence sont par exemple :

- une coopération douteuse de personnes lors d'une enquête;
- une ignorance ou une absence des réponses;
- une interruption de la prise des données (décès d'un patient au cours d'un suivi médical);
- une défaillance dans le système de prise de mesure (défaillance technique temporaire ou non);
- une raisons d'ordre conceptuel (coût d'enquête, de réalisation de prise de mesure);
- une perte de données lors de l'encodage ou du pré-traitement de celles-ci (illisibilité des formulaires d'enquête);
- ...

Dans le contexte d'analyse avec données manquantes, on peut distinguer :

- **le processus de mesure** : le mécanisme gouvernant la distribution des données enregistrées et,
- **le processus de carences** (*missingness*) : le mécanisme gouvernant la distribution des données manquantes.

Comme le souligne [MOLENBERGHS 96b], de la dépendance entre ces deux processus dépend le type d'analyse réalisable. Généralement les deux processus dépendent l'un de l'autre. On peut cependant distinguer, théoriquement, trois types de carences :



1) *Missing Completely At Random* (MCAR) : les deux processus sont statistiquement indépendants (exemple typique : les erreurs techniques). Cette catégorie est la plus simple à traiter. L'apparition d'une donnée manquante dans le tableau de mesures est purement aléatoire.

2) *Missing At Random* (MAR) : le processus de carence dépend des valeurs observées mais pas des valeurs non observées (exemple : une échelle mesure la hauteur d'eau d'une rivière; lors d'une crue l'appareil de mesure est emporté par la force du courant). Une technique intéressante et performante est l'utilisation de l'algorithme *Expectation-Maximisation* (E-M).

3) *Non-Ignorable missingness* (NI) : le processus de carence dépend des valeurs non-observées. Ce type de carence est assez fréquent. Le traitement demande l'emploi de techniques très poussées.

### 3. Les méthodes simples

Avant de développer plus en profondeur l'*E-M algorithm*, attardons-nous sur une série de méthodes simples à utiliser et implémenter.

#### 3.1. Supprimer les données

Une méthode simple, peut-être simpliste, consiste à supprimer les individus pour lesquels les données manquent. Ceci ne peut être réalisé que si le nombre d'individus après suppression reste suffisant. Le risque de perdre de la puissance d'expression existe puisque la taille de l'échantillon diminue. Si la taille de l'échantillon reste suffisamment grande, on peut faire l'hypothèse que non : la suppression des individus pour lesquels les données manquaient n'altère pas l'analyse.

#### 3.2. L'emploi de la moyenne, de la médiane

On peut remplacer les données manquantes par différentes valeurs : la moyenne ou la médiane calculée sur les données présentes au sein d'une même variable.

- La **moyenne** : dans le cas où il y a beaucoup de données, cette approximation semble être assez efficace, rapide, simple. La distribution des données n'est pas trop altérée si les valeurs extrêmes sont peu présentes.

- La **médiane** : les propriétés citées ci-avant pour la moyenne restent valables, mais sans être influencée par les valeurs extrêmes.

Ces deux méthodes de remplacement ne tiennent pas compte de la variabilité de la variable et donc de la distribution, ni du fait que la variable peut être expliquée par d'autres variables de l'étude. Puisqu'il s'agit de remplacer les données manquantes par quelque chose, il est plus intéressant de les remplacer par la valeur la plus probable. Pourquoi ne pas s'aider, se servir des autres variables de l'étude ?



## 4. La régression

On peut remplacer les données manquantes par une valeur estimée à partir des autres variables « corrélées » de l'étude.

On utilise la variable pour laquelle on cherche à remplacer les données manquantes comme variable dépendante et toutes les autres variables (ou quelques-unes judicieusement choisies) comme variables explicatives.

On utilise les observations présentes de la variable en question, on construit un modèle de régression simple ou multiple. On choisit bien évidemment le meilleur modèle parmi l'ensemble des modèles disponibles pour alors prédire les données manquantes.

La régression multiple a pour objectif de caractériser et de formuler la relation qui existe entre plusieurs variables indépendantes et la variable dépendante. Ayant construit un modèle adéquat à partir des observations présentes dans le tableau de nombres de départ, on peut prédire une valeur manquante en appliquant le modèle. En effet, si le modèle est donné par  $Y = \beta_1 X_1 + \dots + \beta_p X_p$  où  $X_1 \dots X_p$  sont les variables explicatives, connaissant les valeurs des variables explicatives pour l'individu  $i$  ( $X_1)_i, \dots, (X_p)_i$  on obtient la valeur manquante par  $(Y)_i = \beta_1(X_1)_i + \beta_2(X_2)_i + \dots + \beta_p(X_p)_i$ .

Les méthodes de régression linéaire ou non linéaire sont largement implémentées au sein de logiciels courants. L'emploi de telles techniques est simple.

Les méthodes utilisant la moyenne, la médiane, et la régression sont raisonnables (en termes de résultats) s'il n'y a pas trop de données manquantes. L'estimation par la moyenne est plus simple que par la médiane, l'utilisation de la régression est la moins aisée des méthodes décrites jusque maintenant. L'ordre de ce classement est inversé lorsqu'il s'agit de qualifier la précision de l'estimation.

L'emploi de la méthode de régression oblige à travailler avec les variables complètes (sans données manquantes pour l'échantillon) comme variables indépendantes. Il peut être judicieux d'employer des modèles intermédiaires pour compléter progressivement le tableau de données.

## 5. Supprimer les variables

L'idée retenue au point 3.1. ci-avant est d'éliminer les individus pour lesquels des données sont manquantes. Si nous appliquons cette suggestion aux variables plutôt que de remplacer une données manquantes, pourquoi ne pas supprimer la variable pour laquelle elle est absente ?

Il ne faut pas supprimer n'importe quelle variable. Pour les sélectionner, on peut, par exemple, appliquer une *Stepwise* (voir chapitre précédent) sur l'ensemble de l'échantillon et l'on relève si la variable a un intérêt explicatif. Si la réponse est négative, cette variable peut être éliminée.



## 6. L'E-M algorithm

Dans le cas où l'ensemble des données manquantes est fortement lié à la valeur de la réponse (par exemple : des employés au salaire élevé refusant de répondre à la question précise sur ce point lors d'une enquête), dans ce cas, les inférences peuvent être fortement biaisées. Actuellement, il n'existe pas de techniques statistiques développées pour ce type de situation.

Par contre, il est possible de traiter des situations pour lesquelles les données manquantes sont aléatoirement distribuées (MAR et MCAR). Car dans ce cas, la distribution des données manquantes n'est pas influencée par les valeurs des variables.

### 6.1. L'algorithme [MOLENBERGHS 96b] et [JOHNSON 92]

Une approche générale de l'estimation de données manquantes a été mise au point par DEMPSTER, LAIRD et RUBIN en 1977. Leur technique appelée l'E-M algorithm est un calcul itératif, comprenant deux étapes :

**l'étape Prediction** : en fonction de  $\theta$ , il s'agit de prédire la contribution d'une donnée manquante aux paramètres statistiques (souvent la moyenne et la variance de la distribution normale);

**l'étape Estimation** : recalculer les paramètres en fonction des valeurs obtenues à l'étape précédente;

avec  $\theta$ , le vecteur paramètre, dans le cas d'une distribution normale, ce paramètre correspond, (c'est l'hypothèse que nous faisons ici pour le développement et l'exemple) à la moyenne et à la variance de la distribution, ou des paramètres qui s'en rapprochent fortement.

On répète le calcul itératif de ces deux étapes jusqu'à ce que les estimations d'une étape ne diffèrent plus de façon significative de celles obtenues à l'itération précédente. L'utilisateur se fixe souvent un seuil d'acceptabilité.

JOHNSON suggère d'utiliser comme paramètres statistiques à estimer :

$$\begin{aligned} T_1 &= \sum_{j=1}^n X_j = n \bar{X} \\ T_2 &= \sum_{j=1}^n X_j X_j' = (n - 1) S + n \bar{X} \bar{X}' \end{aligned}$$

Nous supposons que la variance  $\sigma$  et la moyenne  $\mu$  de la population sont inconnues et doivent être estimées. Une première estimation de ces paramètres est réalisée sur les données non manquantes.



Etape *Prediction* : pour chaque vecteur variable  $x_j$  contenant ces valeurs manquantes,  $x_j = \begin{bmatrix} x_{1j} \\ x_{2j} \end{bmatrix}$ , où  $x_{1j}$  est une donnée manquante et  $x_{2j}$  une donnée valide. En fonction de  $\tilde{\sigma}$  et  $\tilde{\mu}$  de l'étape d'estimation précédente et de la moyenne conditionnelle de la distribution normale de  $x_1$  connaissant  $x_2$ , on estime la contribution de  $x_{1j}$  à  $T_1$  :

$$\tilde{x}_{1j} = E(X_{1j} \mid x_{2j}, \tilde{\mu}, \tilde{\sigma}) = \tilde{\mu}_1 + \tilde{\Sigma}_{12}(\tilde{\Sigma}_{22})^{-1}(\tilde{x}_{2j} - \tilde{\mu}_2) \quad (1)$$

Ensuite, la contribution de  $x_{1j}$  à  $T_2$  est calculée grâce à :

$$Est(x_{1j}(x_{1j})') = E(X_{1j}(X_{1j})' \mid x_{2j}, \tilde{\mu}, \tilde{\sigma}) = \tilde{\Sigma}_{11} - \tilde{\Sigma}_{12}(\tilde{\Sigma}_{22})^{-1}\tilde{\Sigma}_{21} + \tilde{x}_{1j}(\tilde{x}_{1j})' \quad (2)$$

et à :

$$Est(x_{1j}(x_{2j})') = E(X_{1j}(X_{2j})' \mid x_{2j}, \tilde{\mu}, \tilde{\sigma}) = \tilde{x}_{1j}(\tilde{x}_{2j})'$$

Où *Est* signifie : l'opération d'estimation équivalente à  $\tilde{X}$ .

La contribution est sommée pour tous les  $x_j$  contenant des données manquantes.

Les résultats sont combinés avec l'échantillon pour estimer  $\tilde{T}_1$  et  $\tilde{T}_2$ .

Etape *Estimation* : calculer les estimations révisées.

$$\tilde{\mu} = \frac{\tilde{T}_1}{n}, \quad \tilde{\Sigma} = \frac{1}{n}\tilde{T}_2 - \tilde{\mu}\tilde{\mu}'$$

Remarques :

- Le principe de l'estimation avec l'algorithme est très simple, tout spécialement lorsque les données sont en nombre.
- L'algorithme converge toujours vers un maximum local.
- Le taux de convergence est linéaire, et donc assez lent.
- Calculer la précision des estimations n'est pas une tâche triviale.



## 6.2. Un exemple [JOHNSON 92]

Pour comprendre le fonctionnement de cet algorithme, nous illustrons une étape de l'algorithme E-M.

Estimation de la moyenne  $\mu$  de la population et de la covariance  $\sigma$  utilisant le set de données incomplet suivant :

$$\mathbf{X} = \begin{bmatrix} - & 0 & 3 \\ 7 & 2 & 6 \\ 5 & 1 & 2 \\ - & - & 5 \end{bmatrix}$$

Dans ce cas-ci  $n$  (le nombre d'individus) = 4,  $p$  (le nombre de variables) = 3 et certaines observations des vecteurs  $x_1$  et  $x_2$  sont manquantes.

Nous calculons les moyennes des échantillons de base pour les données valables :

$$\tilde{\mu}_1 = \frac{7+5}{2} = 6, \quad \tilde{\mu}_2 = \frac{0+2+1}{3} = 1, \quad \tilde{\mu}_3 = \frac{3+6+2+5}{4} = 4.$$

En substituant ces moyennes aux données manquantes, par exemple  $\tilde{X}_{11} = 6$ , nous pouvons calculer les covariances initiales estimées :

$$\tilde{\sigma}_{11} = \frac{1}{2}; \quad \tilde{\sigma}_{22} = \frac{1}{2}; \quad \tilde{\sigma}_{33} = \frac{5}{2}; \quad \tilde{\sigma}_{12} = \frac{1}{4}; \quad \tilde{\sigma}_{23} = \frac{3}{4}; \quad \tilde{\sigma}_{13} = 1$$

L'étape de prédiction consiste à utiliser les valeur estimées de la moyenne et de la variance pour prédire la contribution des données manquantes aux paramètres  $\mathbf{T}_1$  et  $\mathbf{T}_2$ .

Avec la formule (1), on calcule (la contribution à  $\mathbf{T}_1$ ) :

$$\tilde{x}_{11} = \tilde{\mu}_1 + \tilde{\Sigma}_{12} \tilde{\Sigma}_{22}^{-1} \begin{bmatrix} \tilde{x}_{21} - \tilde{\mu}_2 \\ \tilde{x}_{31} - \tilde{\mu}_3 \end{bmatrix} = 6 + \begin{bmatrix} \frac{1}{4} & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{3}{4} \\ \frac{3}{4} & \frac{5}{2} \end{bmatrix}^{-1} \begin{bmatrix} 0-1 \\ 3-4 \end{bmatrix} = 5,73$$

Avec la formule (2), on calcule (la contribution à  $\mathbf{T}_2$ ) :

$$Est(x_{11}^2) = \tilde{\sigma}_{11} - \tilde{\Sigma}_{12} \tilde{\Sigma}_{22}^{-1} \tilde{\Sigma}_{21} + \tilde{x}_{11}^2 = \frac{1}{2} - \begin{bmatrix} \frac{1}{4} & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{2} & \frac{3}{4} \\ \frac{3}{4} & \frac{5}{2} \end{bmatrix}^{-1} \begin{bmatrix} \frac{1}{4} \\ 1 \end{bmatrix} + (5,73)^2 = 32,09$$



On réalise les mêmes opérations pour les données manquantes de l'individu  $x_4$ .

En utilisant les estimations obtenues ci-dessus, on calcule  $T_1$  et  $T_2$  :

$$T_1 = \sum_{j=1}^n X_j = n \bar{X} = \begin{bmatrix} 24,13 \\ 4,30 \\ 16 \end{bmatrix}$$

$$T_2 = \sum_{j=1}^n X_j X'_j = (n - 1) S + n \bar{X} \bar{X}' = \begin{bmatrix} 148,05 & 27,27 & 101,18 \\ 27,27 & 6,97 & 20,50 \\ 101,18 & 20,50 & 74 \end{bmatrix}$$

Ces calculs clôturent l'étape de prédiction.

L'étape d'estimation fournit les estimations révisées suivantes :

$$\tilde{\mu} = \frac{\tilde{T}_1}{n} = \begin{bmatrix} 6,03 \\ 1,08 \\ 4 \end{bmatrix}$$

$$\tilde{\Sigma} = \frac{1}{n} \tilde{T}_2 - \tilde{\mu} \tilde{\mu}' = \begin{bmatrix} 0,65 & 0,31 & 1,18 \\ 0,31 & 0,58 & 0,81 \\ 1,18 & 0,81 & 2,50 \end{bmatrix}$$

Ceci clôture les calculs de la phase d'estimation.

L'itération suivante des deux étapes continue en utilisant ces nouvelles valeurs. L'algorithme se termine lorsque, d'une itération à la suivante, les éléments  $\tilde{\mu}$  et  $\tilde{\Sigma}$  restent pratiquement inchangés.

Ces calculs sont facilement pris en charge par un ordinateur (heureusement).

Rappelons que l'algorithme est développé en supposant que les données manquantes sont aléatoirement apparues lors de la confection du tableau de données de départ.

## 7. Conclusions

L'analyse un tableau de données incomplet est fréquente. Plusieurs méthodes d'estimation de ces données manquantes sont disponibles. Nous présentons à la figure 9.1. une classification des méthodes présentées dans le cadre de ce chapitre du mémoire.

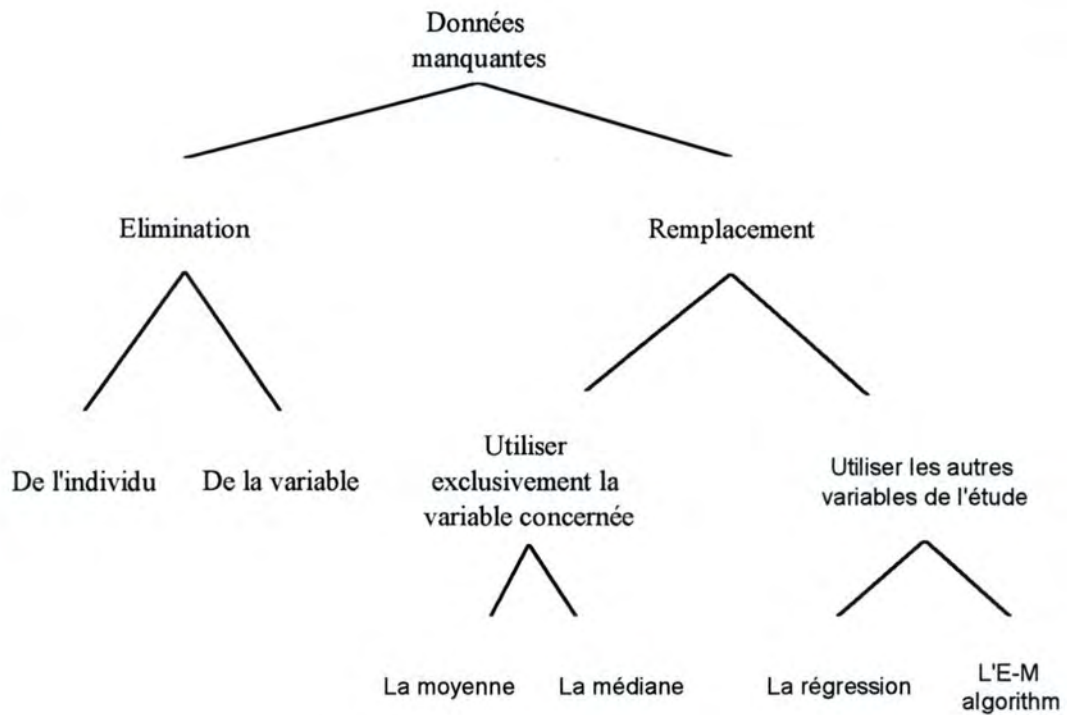


Figure 9.1. : classification des méthodes de résolution du problème des données manquantes.



## CHAPITRE 10 : CONCLUSION

Un objet symbolique est une description qui s'exprime à l'aide d'une conjonction événements (propriétés) portant sur des valeurs prises par des variables.

La théorie des objets symboliques introduite par DIDAY permet d'étendre la problématique, les méthodes et les algorithmes de l'analyse classique des données à des données plus riches. Celles-ci expriment un niveau de connaissance plus élevé que celui fourni par de simples observations exprimées par un vecteur de valeurs quantitatives ou qualitatives.

Il s'agit de donner la possibilité d'utiliser en entrée, des données et des connaissances exprimées par des objets symboliques sans craindre de sortir du carcan tabulaire à  $n$  lignes (les individus) et  $p$  colonnes (les variables) et en évitant de perdre l'information par des modélisations ou codages arbitraires. On s'efforce d'exprimer les résultats sous forme d'objets symboliques, expressions obtenues automatiquement et possédant elles-mêmes un grand pouvoir explicatif.

L'analyse des données symboliques trouve sa place dans un créneau entre l'intelligence artificielle et la statistique, entre les approches logiques, symboliques et numériques. Elle ouvre la voie à un grand champ d'applications : celui du traitement des objets complexes en tenant compte de connaissances non nécessairement d'ordre purement numérique.

Les applications ayant donné des résultats intéressants touchent des domaines aussi divers que les signaux radars, les stratégies de pêche, des scénarii d'accident, des maladies ou des espèces de fleurs. Les deux exemples que nous avons développés permettent de lever un coin du voile sur les potentialités offertes par l'analyse symbolique des données symboliques vis à vis de l'analyse classique de données numériques.

L'analyse des données symboliques a permis de représenter, confirmer, compléter et organiser les scénarii d'accidents issus d'une base de données; des chercheurs ont pu différencier les types d'accidents et spécifier les connaissances apportées par les scénarii.

La seconde analyse résumée rapproche l'analyse des données de l'Intelligence Artificielle. La méthodologie présentée par PERINEL permet de montrer que les deux approches (numérique et symbolique) sont intégrées de façon conjointe et complémentaire :

- L'approche numérique est caractérisée par son efficacité, les méthodes de classification sont utilisées dans la première phase pour déterminer une partition en classes homogènes des pêcheurs.
- L'approche symbolique vient combler certaines lacunes du numérique pur en apportant du sens, de l'explicatif en utilisant les connaissances d'un domaine. Dans le cadre de l'étude, l'apport s'est révélé intéressant pour formaliser de nouveaux concepts que sont les tactiques de pêche exprimées sous forme d'objets explicites (utilisation de la logique modale).



L'objet symbolique est un concept qui permet de décrire de façon assez complète des individus complexes. Que ce soit sous forme de tableaux ou d'assertions, percevoir l'information véhiculée n'est pas nécessairement chose aisée : une représentation graphique s'impose.

La complexité de l'objet symbolique est telle qu'il est nécessaire d'opter pour une représentation propre. L'idéal est de représenter un maximum d'informations sans surcharges.

Sur le marché, il n'existe pas de logiciels permettant ce type de représentation (multivariée, multi-échelle, ...). Dans le cadre de leurs recherches, M. NOIRHOMME-FRAITURE et M. ROUARD ont développé une solution pour la représentation d'objets statistiques complexes : « l'étoile zoom ».

Cette représentation en 'Etoile Zoom' donne une image synthétique d'objets multivariés complexes. Ces graphiques sont des étoiles compte tenu de la représentation d'axes radiaires.

Elles sont qualifiées par « zoom » parce que, en fonction de ses besoins et de ses desiderata, l'utilisateur peut en sélectionnant un axe particulier obtenir des informations plus précises sur la variable en question : dépendance entre variables, représentation supplémentaire des histogrammes de poids.

Le son est utilisé pour présenter de l'information qui ne pourrait être l'objet de représentation via un mode de visualisation classique ou d'informations difficiles à discerner, telles que les données numériques multidimensionnelles. Le son est habituellement un complément aux *outputs* visuels parce qu'il augmente la quantité d'informations communiquées aux utilisateurs et/ou réduit la quantité d'informations que l'utilisateur doit traiter en mode visuel.

Nous avons divisé le chapitre selon le type d'informations véhiculé par le son. Utilisé soit comme moyen de signalisation, d'alerte, soit comme moyen de compréhension de phénomènes, le son aide à interpréter une figure, une représentation.

La combinaison des informations graphiques et sonores sur une même interface est devenue naturelle. Chaque jour, les deux sens sont associés pour permettre la perception optimale d'informations complémentaires du monde. Les deux sens sont interdépendants.

Le système visuel nous donne des détails d'un foyer, là où le système auditif fournit des informations générales de ce qui nous entoure, nous alertant de choses que nous ne pouvons voir. Les deux sens combinés nous fournissent toute l'information (parfois plus) dont nous avons besoin pour comprendre l'environnement.

Au-delà de l'aide à la perception que sont les sons de signalisation, ne peut-on pas aller plus loin ? Comment utiliser le son comme moyen de perception, de représentation et pas seulement comme moyen de signalisation ?

Les *feed-back* audios ont prouvé leurs apports significatifs pour des interfaces homme-machine d'applications interactives. L'emploi du son est justifié lorsque la présentation de l'information ne peut être formulée autrement, ou lorsque l'attention de l'utilisateur doit être focalisée sur un sujet particulier, ou encore si celui-ci effectue une autre tâche en même temps.



Nous pouvons retenir de ces expériences que :

- les *earcons* facilitent la communication d'informations;
- l'utilisation de séquences de notes est préférable aux sons simples (simplistes);
- les capacités de mémorisation sont d'autant plus importantes que les séquences sont différentes entre elles;
- le fait d'être musicien n'augmente pas de façon significative les performances d'utilisation d'interfaces sonores.

Nous synthétisons pour chaque paramètre accessible à l'utilisateur (timbre, note, registre, rythme, intensité) les principales recommandations d'emploi. Il s'agit de conseils de niveau conceptuel. Du respect de ceux-ci dépendent la réussite et l'efficacité de l'interface. Le respect de ces règles n'est pas strict et absolu.

L'utilisation plus ou moins harmonieuse des notes musicales, plutôt que de salves de bruits, est plus performante. Réaliser des séquences de notes, tout en respectant les lignes directrices citées procure une réelle augmentation de performance. Les auteurs soulignent qu'un temps d'apprentissage est nécessaire, mais c'est un investissement rentable.

Les développeurs peuvent créer des interfaces utilisant le son, c'est un bon moyen de communication. L'interprétation de figures 3D peut être simplifiée par l'emploi d'environnements sonores structurés.

L'emploi du son au sein des interfaces homme-machine représente un véritable potentiel d'aide à la représentation d'objets complexes, comme les objets symboliques. Comment le son peut-il aider l'utilisateur dans sa tâche d'analyse, d'interprétation, de compréhension des étoiles zoom ?

L'idée d'un son typique pré-programmé est intéressante pour la mise automatique en évidence de certains événements (événements de l'environnement de travail et événements propres aux étoiles zoom). Les *earcons* font référence à des bruits familiers, le but est d'associer une action, un événement à un « environnement » sonore connu de l'utilisateur. D'autres événements, d'autres actions peuvent faire l'objet d'adjonctions sonores. L'utilisateur doit pouvoir, au moins, supprimer la génération sonore de ces séquences et, éventuellement les modifier ou en sélectionner d'autres en remplacement.

Le son guide l'utilisateur par le son dans sa localisation sur la représentation de l'étoile, son interprétation de la forme de l'étoile. Les déplacements de l'utilisateur selon l'axe horizontal sont détectés par une variation de la balance et ses mouvements verticaux par l'intensité d'émission. Les variables pour lesquelles une pondération est utilisée sont représentées partiellement. Un axe en pointillé signale visuellement la caractéristique mais la représentation de l'étoile ne tient compte que de la valeur munie du poids le plus important. Le son pallie à cette carence. Un environnement musical aide l'utilisateur dans son interprétation de la représentation de la 3D des poids des variables. Nous décomposons en tranches égales les valeurs de poids potentielles souvent exprimées en pour-cent. Nous associons à chaque tranche un instrument différent.

Nous avons montré l'apport du son pour aider l'utilisateur dans sa tâche d'analyse, d'interprétation, de compréhension des étoiles zoom.



Depuis longtemps les couleurs sont utilisées au sein des applications interactives. Le matériel informatique actuel permet l'emploi de couleurs dans toutes les tâches habituelles. Les possibilités offertes par des périphériques de moins en moins onéreux autorisent la circulation de documents couleurs et un échange direct avec l'ordinateur. A l'inverse du son, l'emploi de la couleur est largement répandu dans toutes les applications et tous les domaines.

Nous avons montré l'utilité, l'apport de la couleur par un exemple expérimental. Nous nous intéressons au « mécanisme » oculaire de perception des couleurs et à certains phénomènes d'illusion qu'il convient d'éviter. En fonction des trois paramètres que sont la texture, l'intensité et les nuances nous mettons en évidence des recommandations d'utilisation de la couleur. Pour les couleurs (bleu, vert, jaune, rouge, noir, blanc), nous avons identifié son apport et son applicabilité et édicté des recommandations d'utilisation. En plus du côté esthétique, l'apport de la couleur dans une application est significatif.

L'apport de la troisième dimension se justifie pour les objectifs suivants :

- mise en évidence des **détails**;
- **transmission** entre utilisateurs de l'information;
- favoriser la **mémorisation** de la représentation.

Lorsque le concepteur intègre la 3D dans les représentations, il convient de :

- préférer les **vues plongeantes**;
- d'employer les **histogrammes** pour les comparaisons (graphiques en barres).

Les situations où des éléments de la représentation 3D sont dissimulés derrière d'autres représentent les limites de la représentation 3D. Dans le cadre des représentations 3D des étoiles zoom ces difficultés sont contournées : la rotation et la modification de l'angle de vue sont retenus.

En marge du sujet principal de ce mémoire, nous nous sommes intéressés aux données manquantes et à la sélection des variables. En effet, l'un des buts de la représentation d'objets symboliques étant de les comparer, l'apparition de « trous » dans le tableau de nombres de base pose un sérieux problème. D'autre part la représentation graphique donne une vue d'ensemble, le nombre de variables représentables sur celui-ci est forcément limité.



## BIBLIOGRAPHIE

### Citée

- [AGRESTI 90] AGRESTI A., *Categorical data analysis*. Univesity of Florida, Wiley Interscience, John Wiley & sons, 1990, 558 pages.
- [BEAUDOUIN-LAFON 96] BEAUDOUIN-LAFON M., CONVERSY S., *Auditory illusions for audio feed-back*. CHI 96 Conference short papers, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 299-300.
- [BREWSTER 93] BREWSTER S. A., WRIGHT P. C., EDWARDS A. D. N., *An evaluation of earcons for use in auditory human-computer interfaces*. INTERCHI 93 Conference proceedings, Conference on Human Factors in Computing Systems, 24-29 avril 1993, pages 222-227.
- [BREWSTER 94] BREWSTER S. A., WRIGHT P. C., EDWARDS A. D. N., *The design and evaluation of an auditory-enhanced scrollbar*. CHI 94 Conference papers, Conference on human factors in computing systems, Boston, Massachusetts, 24-28 avril 1994, pages 173-179.
- [BREWSTER 95] BREWSTER S.A. WRIGHT P. C., EDWARDS A. D. N., DIX A. J., *The sonic enhancemant of graphical buttons*. Interact '95 Conference papers, Human-Computer Interaction, Chapman & Hall, Londres, 1995, pages 43-48.
- [BREWSTER 97] BREWSTER S.A. MURRAY G. *Making memus musical*. CHI 97 Conference Proceedings, "Human Factors in Computer Systems", Atlanta, 22-27 mars 1997, pages 389-396.
- [COWPERTHWAIT 96] COWPERTHWAIT D.J., SHEELAGH M., CARPENDALE T., FRACCHIA F., *Visual access for 3D data*. CHI 96 Conference short papers, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 175-176.
- [DE CARVALHO 94] DE CARVALHO F.A.T., *Un indice de dissimilarité entre objets symboliques booléens basé sur l'extension*. Actes des 4<sup>èmes</sup> Journées sur l'Induction Symbolique/Numérique, E. Diday, Y. Kodratoff et M. Moulet (éditeurs), Orsay, France, 14-15 mars 1994, pages 29-44.



- [DENIS 92] DENIS J., *Géographie de la Belgique*. Ouvrage collectif sous la direction de Jacques Denis. Crédit Communal, Bruxelles, 1992, 622 pages.
- [DIDAY 87] DIDAY E., KODRATOFF Y., *Introduction à l'approche symbolique en Analyse des Données*. Actes des journées symboliques - numériques pour l'apprentissage à partir de données. CEREMADE. Université Paris IX Dauphine, 1987, 327 pages.
- [DIDAY 91] DIDAY E., KODRATOFF Y., *Induction symbolique et numérique à partir de données*. Volume I. Cépaduès-Editions, 1991 460 pages.
- [DIDAY 93] DIDAY E., *Quelques aspects de l'analyse des données symboliques*. Rapport de recherche n° 1937. Institut National de Recherche en Informatique et en Automatique, 1993, 19 pages.
- [DIGIANO 93] DIGIANO C., BAECKER R., OWEN N., *LogoMedia : A sound-enhanced programming environment for monitoring program behaviour.*, INTERCHI 93 Conference proceedings, Conference on Human Factors in Computing Systems, 24-29 avril 1993, pages 301-302.
- [DOUGLAS 96] DOUGLAS S., KIRKPATRICK T., *Do color models really make a difference ?* CHI 96 Conference papers, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 399-405.
- [ENKLAAR 95] ENKLAAR C., *Brightness/color illusions*. Tutorial disponible sur le réseau à l'adresse suivante : <http://hyperg.uni-paderborn.de/>, septembre 1995, 4 pages.
- [GAVER 93] GAVER W. W., *Synthesizing auditory icons*. INTERCHI 93 Conference proceedings, Conference on Human Factors in Computing Systems, 24-29 avril 1993, pages 228-235.
- [JACKSON 94] JACKSON J. A., FRANCONI J. M., *Synchronization of visual and aural parallel program performance data*. Auditoriy display; Kramer (Editor), Addison-Wesley, Reading MA, 1994, pages 369-416.
- [JOHNSON 92] JOHNSON R. A., WICHERN D. W., *Applied multivariate statistical analysis*. Third edition, Prentice Hall International Editions, 1992, 642 pages.
- [LÊ 96] LÊ T., *Interface entre un SGBD Oracle et un nouveau logiciel d'analyse de données*. Rapport de stage MIAIF 2<sup>ème</sup> année, 1996, 55 pages.
- [LEIMANN 95] LEIMANN E., SCHULZE H.H., *Earcons and icons : an experimental study*. Interact '95 Conference papers, Human-Computer Interaction, Chapman & Hall, Londres, 1995, pages 49-54.



- [LEVY 96] LEVY E., ZACKS J., TVERSKY B., SCHIANO D., *Gratuitous graphics ? Putting preferences in perspective*. CHI 96 Conference papers, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 42-49.
- [LYNCH 94a] LYNCH P.J., *Visual design for the user interface. Part 1: design fundamentals*. Journal of biocommunications n°21 (1), 1994, pages 22-30.
- [LYNCH 94b] LYNCH P.J., *Visual design for the user interface. Part 2: graphics in the interface*. Journal of biocommunications n°21 (2), 1994, pages 6-15.
- [MACCHI 96] MACCHI W. S., *Colors models*. Tutorial disponible sur le réseau à l'adresse suivante : <http://longwood.cs.ucf.edu/~macchi/research>, 1996, 5 pages.
- [MEREU 96] MEREU S.W., KAZMAN R., *Audio enhanced 3D interfaces for visually impaired users*. CHI 96 Conference papers, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 72-78.
- [MOLENBERGHS 96a] MOLENBERGHS G., *Multivariate analysis : Introduction to multivariate data analysis*, Part 1, Faculteit Wetenschappen, Limburgs Universitair Centrum, 1996, 189 pages.
- [MOLENBERGHS 96b] MOLENBERGHS G., *Multivariate analysis : Introduction to multivariate data analysis*, Part 2, Faculteit Wetenschappen, Limburgs Universitair Centrum, 1996, 348 pages.
- [MULLET 96] MULLET K.E., *Designing visual interfaces: how to create communication-oriented solutions*. CHI 96 Conference tutorials, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 332-333.
- [MURAKAMI 94] MURAKAMI T., NAKAJIMA N., *Direct and intuitive input device for 3D shape deformation*. CHI 94 Conference short papers, Conference on human factors in computing systems, Boston, Massachusetts, 24-28 avril 1994, pages 465-470.
- [MYNATT 94] MYNATT E. D., *Designing with auditory icons : how well do we identify auditory cues ?*. CHI 94 Conference short papers, Conference on human factors in computing systems, Boston, Massachusetts, 24-28 avril 1994, pages 269-270.



- [NETER 96] NETER J., KUTNER M. H., NACHTSHEIM C. J., WASSERMAN W., *Applied linear statistical models*. Fourth edition, IRWIN, 1996, 1408 pages.
- [NOIRHOMME-FRAITURE 96] NOIRHOMME-FRAITURE M., ROUARD M., *L'étoile zoom: une solution pour la représentation d'objets statistiques complexes*. RP-96-038 disponible sur le réseau à l'adresse suivante : <http://www.info.fundp.ac.be/cgi-bin/pub-RP>, Institut d'Informatique, FUNDP Namur, Décembre 1996, 6 pages.
- [NOIRHOMME-FRAITURE 97a] NOIRHOMME-FRAITURE M., ROUARD M., *Zoom star : a solution to complex statistical object representation*. RP-97-005 disponible sur le réseau à l'adresse suivante : <http://www.info.fundp.ac.be/cgi-bin/pub-RP>, Institut d'Informatique, FUNDP Namur, 1997, 2 pages.
- [NOIRHOMME-FRAITURE 97b] NOIRHOMME-FRAITURE M., ROUARD M., *Computer graphics for symbolic objects*. RP-97-006 disponible sur le réseau à l'adresse suivante : <http://www.info.fundp.ac.be/cgi-bin/pub-RP>, Institut d'Informatique, FUNDP Namur, 1997, 4 pages.
- [PERINEL 92] PERINEL E., *Analyse numérique / symbolique des tactiques de pêche artisanale au Sénégal*. Rapport de stage de DEA de Mathématique Appliquées aux Sciences Economiques. Université Paris IX Dauphine, 1992, 74 pages.
- [PISSART 92] PISSART A., JUVIGNE E., *Méthodes d'étude des formations détritiques continentales*. Cours de spécialisation en géographie physique. Université de Liège, 1992-1993, 142 pages.
- [REGNIER 92] REGNIER A., *Analyse Symbolique de Scénarios d'Accidents*. Rapport de stage de DEA MAI option Intelligence Artificielle. Université Paris IX Dauphine, 1992, 85 pages.
- [RIGAS 97] RIGAS D. I., ALTY J. L., *The use of music in a graphical interface for the visually impaired*. CHI 97 Conference Proceedings, "Human Factors in Computer Systems", Atlanta, 22-27 mars 1997, pages 228-235.
- [ROSE 96] ROSE J.J., *Perception theories*. Tutorial disponible sur le réseau à l'adresse suivante : <http://hyperg.uni-paderborn.de/>, septembre 1996, 7 pages.
- [SAKAI 96] SAKAI A., *Flying fingers : a tool for three-dimensional shared workspace*. CHI 96 Conference short papers, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 295-296.



- [SHERRY 95] SHERRY, L. Perception Principles. Tutorial disponible sur le réseau à l'adresse suivante : <http://www.cudenver.edu/~lsherry/perception.html>, septembre 1995, 10 pages.
- [TAPP 92] TAPP R., KAZMAN R., *Determining the usefulness of colour and fonts in programming task*. Article disponible sur le réseau à l'adresse suivante : <ftp://cs-archive.uwaterloo.ca/cs-archive/Index>, 1995, 8 pages.
- [UNRUH 96] UNRUH P., *Perception theories*. Tutorial disponible sur le réseau à l'adresse suivante : <http://hyperg.uni-paderborn.de/>, juillet 1996, 12 pages.
- [VANDERDONCKT 94] VANDERDONCKT J., *Guide ergonomique des interfaces homme-machine Guide pratique pour la conception ergonomique des interfaces homme-machine pour les applications interactives*. Collection « Travaux de l'Institut d'Informatique », n°13, FUNDP, Namur, 1994, 625 pages.
- [VANDERDONCKT 96a] VANDERDONCKT J., BODART F., *The « Corpus Ergonomicus » : a comprehensive and unique source for human-machine interface guidelines*. RP-96-022 disponible à l'adresse suivante : <http://www.info.fundp.ac.be/cgi-bin/pub-RP>, Institut d'Informatique, FUNDP Namur, Juillet 1996, 6 pages.
- [VANDERDONCKT 96b] VANDERDONCKT J., *Cours d'Interface homme-machine de 2<sup>ème</sup> Licence*, Institut d'Informatique, FUNDP Namur, 1995-1996.
- [VENOLIA 93] VENOLIA D., *Facile 3D direct manipulation*. INTERCHI 93 Conference proceedings, Conference on Human Factors in Computing Systems, 24-29 avril 1993, pages 31-36.

## Compulsée

- [BRITO 93] BRIT P., CHABANON C., *Analyse de données numérique et symbolique appliquée au confort automobile*. INRIA et Renault S.A Centre technique d'Aubevoye, 1993, 20 pages.
- [CHIDANANDA 91] CHIDANANDA K., DIDAY E., *Unsupervised learning through symbolic clustering*. Pattern recognition letters 12, North Holland, 1991, pages 259-264.
- [CHIDANANDA 92] CHIDANANDA K., DIDAY E., *Symbolic clustering using a new similarity measure*. IEEE transactions on systems, man, and cybernetics, vol. 22, No 2, ars/avril 1992, pages 368-378.



- [GAVER 93] GAVER W.W., *Synthesizing auditory icons*. INTERCHI 93 Conference proceedings, Conference on Human Factors in Computing Systems, 24-29 avril 1993, pages 228-235.
- [KAZMAN 91] KAZMAN R., *Babel : a psychologically plausible cross-linguistic model of lexical and syntactic acquisition*. L.A. Birnbaum and G.C. Collins, éditeurs, Proceedings of the eighth International workshop, San Mateo, Californie, 1991, pages 75-79.
- [KAZMAN 94] KAZMAN R., *Simulating the child's acquisition of the lexicon and syntax-experiences with Babel*. Machine learning. Kluwer academic publishers, Boston, 1994, pages 1-36.
- [KAZMAN 96] KAZMAN R., *Accessing multimedia through concept clustering*. Article disponible à l'adresse suivante : <http://www.sei.cmu.edu/staff/rkazman/>. University de Wterloo, Ontario, Canada, 1996, 8 pages.
- [LENNART 93] LENNART E. F., BROWN C. G., STALH O., CARLSSON C., *A space based model for user interaction in shared synthetic environments*. INTERCHI 93 Conference proceedings, Conference on Human Factors in Computing Systems, 24-29 avril 1993, pages 43-48.
- [LEVENTHAL 95] LEVENTHAL L., TEASLEY B., INSTONE K., STONE D., *Searching without a keyboard in a multimedia environment*. Interact '95 Conference papers, Human-Computer Interaction, Chapman & Hall, Londres, 1995, pages 241-246.
- [MYNATT 94a] MYNATT E. D., WEBER G. *Nonvisual presentation of graphical user interfaces : contrasting two approaches*. CHI 94 Conference papers, Conference on human factors in computing systems, Boston, Massachusetts, 24-28 avril 1994, pages 166-172.
- [NOIRHOMME-FRAITURE 95] NOIRHOMME-FRAITURE M., *Cours de Statistiques Appliquées de 2<sup>ème</sup> Licence*, Institut d'Informatique, FUNDP Namur, 1994-1995.
- [NOIRHOMME-FRAITURE 96b] NOIRHOMME-FRAITURE M., GOFFINET L., *Automatic hypertexte link generation based on similarity measures between documents*. RP-96-034 disponible sur le réseau à l'adresse suivante : <http://www.info.fundp.ac.be/cgi-bin/pub-RP>, Institut d'Informatique, FUNDP Namur, Décembre 1996, 7 pages.
- [RAMSTEIN 95] RAMSTEIN C., *A multipodal user interface system with force feedback and physical models*. Interact '95 Conference papers, Human-Computer Interaction, Chapman & Hall, Londres, 1995, pages 157-168.



- [SCHUMANN 96] SCHUMANN J., STROTHOTTE T., LASER S., *Assessing the effect of non-photorealistic rendered images in CAD*. CHI 96 Conference papers, Conference on human factors in computing systems, Vancouver, 13-18 avril 1996, pages 35-41.
- [SIKORA 97] SIKORA C. A., ROBERTS L.A. *Defining a family of feed-back signals for multimedia communication devices*. CHI 97 Conference Proceedings, "Human Factors in Computer Systems", Atlanta, 22-27 mars 1997, pages 373-380.
- [TILO 96] TILO, F., *The Perception theories*. Tutorial disponible sur le réseau à l'adresse suivante : <http://hyperg.uni-paderborn.de/>, juillet 1996, 3 pages.
- [UNRUH 96b] UNRUH P., *Marr's computational approach*. Tutorial disponible sur le réseau à l'adresse suivante : <http://hyperg.uni-paderborn.de/>, juillet 1996, 4 pages.
- [ZHAI 94] ZHAI S., BUXTON W., MILGRAM P., *The « silk cursor » : investigating transparency for 3D target acquisition*. CHI 94 Conference papers, Conference on human factors in computing systems, Boston, Massachusetts, 24-28 avril 1994, pages 459-464.



## **TABLE DES MATIERES**

### **Remerciements**

### **Introduction** 01

#### **CHAPITRE 01 :**

1. Introduction	05
2. Les concepts de la théorie des objets symboliques	05
2.1. Les variables	05
2.2. Les objets symboliques	06
2.3. Les événements élémentaires	07
2.4. Les objets assertion	07
2.5. Les objets hordes	08
2.6. Les objets de synthèse	09
2.7. Les objets symboliques munis de méthodes et de propriétés	09
2.8. La connaissance supplémentaire	10
2.9. Les traitements des données symboliques	11
2.10. Propriétés des objets symboliques	12
2.11. Qualité des objets symboliques	13
2.12. Qualité des classes d'objets symboliques	15
2.14. Une extension à des objets symboliques modaux	16
3. Conclusion	21

#### **CHAPITRE 02 : L'ANALYSE DES OBJETS SYMBOLIQUES**

1. Introduction	23
2. L'approche symbolique : l'analyse des données	23
2.1. Les quatre type d'analyse de données	23
2.2. L'analyse des données symboliques par rapport à d'autres disciplines	24
2.3. Les six étapes d'une analyse des données symboliques [DIDAY 93]	24
3. L'approche symbolique : Exemples d'application	25
3.1. Scénarii d'accidents : un outil pour les diagnostics de sécurité	25
3.2. Analyse symbolique des tactiques de pêche artisanale au Sénégal	27
4. Dissimilarité	30
4.1. Le potentiel de description d'un objet assertion booléen	31
4.2. Le calcul de la proximité entre objets assertion booléens	32
5. Conclusions [DIDAY 93]	35

#### **CHAPITRE 03 : LA REPRESENTATION EN ETOILE**

1. Introduction	37
2. L'Etoile Zoom	37
2.1. Introduction à l'Etoile Zoom	38
2.2. La représentation des objets symboliques	39
2.3. Plus avant dans la représentation	41
2.4. L'implémentation	43
3. Conclusions	44



## **CHAPITRE 04 : L'INTERFACAGE**

1. Introduction	45
2. Concevoir une interface	45
2.1. Les règles de conception	46
3. Manipuler les objets d'une interface	48
4. Interpréter les éléments d'une interface	50
4.1. La perception des formes	50
4.2. La perception des images	51
4.3. Le contraste	51
4.4. La perception des graphiques, diagrammes	52
4.5. La perception des sons	52
5. Conclusions	52

## **CHAPITRE 05 : L'UTILISATION DU SON**

1. Introduction	53
2. Le son : moyen de signalisation	54
2.1. Introduction	54
2.2. Les expériences	54
2.3. Conclusion	57
3. Le son : moyen de communication d'informations	58
3.1. Introduction	58
3.2. Les expériences	58
3.3. Synthèse des résultats	68
4. Conclusions	68

## **CHAPITRE 06 : L'UTILISATION DES COULEURS**

1. Introduction	75
2. La perception des couleurs	76
2.1. La théorie de YOUNG-HELMHOLTZ	77
2.2. La théorie de HERING	77
3. Les illusions	78
4. Les couleurs et les étoiles zoom	79
4.1. La texture	79
4.2. L'intensité	79
4.3. Les nuances	80
5. Conclusion	81

## **CHAPITRE 07 : L'UTILISATION DU RELIEF**

1. Introduction	83
2. Les circonstances d'emploi de la 3D	83
2.1. Introduction	83
2.2. Les expériences	83
2.3. Synthèse	84
3. 3D et visibilité	84
4. Les manipulations	86
5. Conclusion	87

## CHAPITRE 08 : LA SELECTION DE VARIABLES

1. Introduction	89
2. Des critères de sélection	89
2.1. Le critère $R_p^2$ ou $SSE_p$	89
2.2. Le critère $MSE_p$ ou $R_a^2$	90
2.3. Le critère $C_p$	91
2.4. Le critère $PRESS_p$ .	91
3. Les procédures automatiques de sélection de variables	92
3.1. Introduction	92
3.2. Forward Stepwise Regression	93
3.3. D'autres procédures automatiques de recherche	94
4. Conclusions	94

## CHAPITRE 09 : LES DONNEES MANQUANTES

1. Introduction	95
2. Les données manquantes	95
3. Les méthodes simples	96
3.1. Supprimer les données	96
3.2. L'emploi de la moyenne, de la médiane	96
4. La régression	97
5. Supprimer les variables	97
6. L'E-M algorithm	98
6.1. L'algorithme [MOLENBERGHS 96b] et [JOHNSON 92]	98
6.2. Un exemple [JOHNSON 92]	100
7. Conclusions	102

<i>Conclusions</i>	103
--------------------	-----

<i>Bibliographie</i>	107
----------------------	-----

<i>Table des matières</i>	115
---------------------------	-----

Annexe



Annexe 01.

Les méthodes que nous présentons dans les chapitres 08 et 09 s’adressent en particulier à des variables quantitatives. Les variables des objets symboliques ne sont pas nécessairement de ce type. Des opérations de codage et de transformation s’imposent pour permettre l’application des méthodes décrites. Nous reprenons ci-dessous les opportunités de modifications. Ces remarques sont valables pour les deux chapitres traitant respectivement de la sélection des variables et des données manquantes.

Pour *les variables quantitatives continues*, la nécessité d’un codage ne se présente pas.

Pour *les variables quantitatives discrètes*, la Stepwise est applicable. Seules les méthodes sur les données manquantes posent problème. En fait il faut fixer en plus des « *cut points* » qui permettent d’attribuer à une valeur prédite, pour une observation manquante, une valeur observable. Par exemple : si les valeurs observables sont 1, 2, 3, 4, ... et si la valeur pour une donnée manquante est de 3,45, celle-ci doit être ramenée à une des valeurs possibles (3 ou 4). Faut-il tronquer ou arrondir les décimales ? La technique doit être déterminée a priori.

Lorsque *la variable quantitative est de type intervalle* : nous proposons par exemple de scinder cette variable en deux sous variables. La première correspond au minimum de l’intervalle et la seconde au maximum.

Exemple :

$[y_2 = [a,b]]$  devient  $[y'_2 = a] \wedge [y''_2 = b]$

Lorsque nous considérons *les variables qualitatives nominales* : pour un individu donné et pour chaque catégorie, nous concevons de coder sa réponse sous forme de 0 et de 1 en donnant la valeur 1 à la catégorie choisie. Chaque variable est scindée en autant de nouvelles variables qu’il y a de catégories dans la variable originale. Par exemple :

Si une enquête est constituée de trois questions, la question un a trois catégories possibles, la deuxième quatre et la dernière deux. Chaque individu est repéré par une ligne de 0 et de 1 :

	Question 1	Question 2	Question 3
Individu <i>i</i>	0 0 1	0 0 1 0	1 0

Remarques :

- Cette technique est déjà utilisée en analyse classique, par exemple pour l’analyse des correspondances multiples, la régression multiple.
- Si plusieurs catégories sont présentes simultanément (ce qui est permis pour les objets symboliques), il faut coder à 1 toutes celles présentes.

Lorsqu'une des *variables est munie d'un poids* pour ses modalités (par exemple : une probabilité), nous proposons de décomposer cette variable et de donner un poids à toutes les variables : le poids lié au mode pour les variables qui en étaient munies et un poids de 1 pour toutes les autres variables de l'analyse.

Exemple :

$[y_2 = \text{âge}] \wedge [y_2 (\text{couleur des yeux}) = \{0,90\{\text{bleu}\}; 0,10\{\text{brun}\}\}]$  devient

$[y_2 = \text{âge}]$  avec un poids de 1  
 $\wedge$   $[y_2 = \{\text{bleu}\}]$  avec un poids de 0,90  
 $\wedge$   $[y_2 = \{\text{brun}\}]$  avec un poids de 0,10

Pour *les variables ordinales*, nous proposons de travailler avec les rangs comme données d'une variable quantitative.

Remarque :

Si une variable est scindée en deux ou plusieurs variables, pour garder la cohérence, il faut forcer la simultanéité d'action des variables résultantes. Par exemple : si  $y_1$  est scindée en  $y'_1$  et  $y''_1$ , lors de l'étape de la procédure *Stepwise*, si  $y'_1$  est sélectionnée pour participer (resp. être enlevée) au modèle, il faut forcer l'intervention (resp. l'élimination) de  $y''_1$ <sup>1</sup>. Les deux variables ( $y'_1$  et  $y''_1$ ) doivent toujours être présentes (ou absentes) en même temps. Le logiciel SAS, par exemple, permet de forcer la présence de variables dans le modèle. Avec le logiciel GLIM<sup>2</sup>, la procédure *Stepwise* s'implémente « interactivement » : l'utilisateur propose son choix de variables, le logiciel calcule les coefficients et permet de déterminer les critères introduits au chapitre 08.

<sup>1</sup> Dans ce cas, ce n'est plus seulement l'apport de  $y'_1$  qui doit être calculé, mais l'apport de  $y'_1$  et  $y''_1$ .

<sup>2</sup> Generalised linear interactive model.